

IEEE Std 1003.1-2001
(Revision of IEEE Std 1003.1-1996
and IEEE Std 1003.2-1992)

Open Group Technical Standard
Base Specifications, Issue 6

1003.1TM

**Standard for Information Technology —
Portable Operating System Interface (POSIX[®])**

Rationale (Informative)

Approved 12 September 2001
The Open Group

IEEE Sponsor

Portable Applications Standards Committee
of the
IEEE Computer Society

Approved 6 December 2001
IEEE-SA Standards Board



THE *Open* GROUP

Abstract

This standard defines a standard operating system interface and environment, including a command interpreter (or “shell”), and common utility programs to support applications portability at the source code level. It is the single common revision to IEEE Std 1003.1-1996, IEEE Std 1003.2-1992, and the Base Specifications of The Open Group Single UNIX[®]† Specification, Version 2. This standard is intended to be used by both applications developers and system implementors and comprises four major components (each in an associated volume):

- General terms, concepts, and interfaces common to all volumes of this standard, including utility conventions and C-language header definitions, are included in the Base Definitions volume.
- Definitions for system service functions and subroutines, language-specific system services for the C programming language, function issues, including portability, error handling, and error recovery, are included in the System Interfaces volume.
- Definitions for a standard source code-level interface to command interpretation services (a “shell”) and common utility programs for application programs are included in the Shell and Utilities volume.
- Extended rationale that did not fit well into the rest of the document structure, containing historical information concerning the contents of this standard and why features were included or discarded by the standard developers, is included in the Rationale (Informative) volume.

The following areas are outside the scope of this standard:

- Graphics interfaces
- Database management system interfaces
- Record I/O considerations
- Object or binary code portability
- System configuration and resource availability

This standard describes the external characteristics and facilities that are of importance to applications developers, rather than the internal construction techniques employed to achieve these capabilities. Special emphasis is placed on those functions and facilities that are needed in a wide variety of commercial applications.

Keywords

application program interface (API), argument, asynchronous, basic regular expression (BRE), batch job, batch system, built-in utility, byte, child, command language interpreter, CPU, extended regular expression (ERE), FIFO, file access control mechanism, input/output (I/O), job control, network, portable operating system interface (POSIX[®]†), parent, shell, stream, string, synchronous, system, thread, X/Open System Interface (XSI)

† See **Trademarks** (on page xxv).

Rationale (Informative)

Published 6 December 2001 by the Institute of Electrical and Electronics Engineers, Inc.
3 Park Avenue, New York, NY 10016-5997, U.S.A.
ISBN: 0-7381-3050-8 PDF 0-7381-3010-9/SS94956 CD-ROM 0-7381-3129-6/SE94956
Printed in the United States of America by the IEEE.

Published 6 December 2001 by The Open Group
Apex Plaza, Forbury Road, Reading, Berkshire RG1 1AX, U.K.
Document Number: C953
ISBN: U.K. 1-85912-262-0 U.S. 1-931624-10-0
Printed in the U.K. by The Open Group.

All rights reserved. No part of this publication may be reproduced in any form, in an electronic retrieval system or otherwise, without prior written permission from both the IEEE and The Open Group.

Portions of this standard are derived with permission from copyrighted material owned by Hewlett-Packard Company, International Business Machines Corporation, Novell Inc., The Open Software Foundation, and Sun Microsystems, Inc.

Permissions

Authorization to photocopy portions of this standard for internal or personal use is granted provided that the appropriate fee is paid to the Copyright Clearance Center or the equivalent body outside of the U.S. Permission to make multiple copies for educational purposes in the U.S. requires agreement and a license fee to be paid to the Copyright Clearance Center.

Beyond these provisions, permission to reproduce all or any part of this standard must be with the consent of both copyright holders and may be subject to a license fee. Both copyright holders will need to be satisfied that the other has granted permission. Requests to the copyright holders should be sent by email to austin-group-permissions@opengroup.org.

Feedback

This standard has been prepared by the Austin Group. Feedback relating to the material contained in this standard may be submitted using the Austin Group web site at <http://www.opengroup.org/austin/defectform.html>.

IEEE

IEEE Standards documents are developed within the IEEE Societies and the Standards Coordinating Committees of the IEEE Standards Association (IEEE-SA) Standards Board. The IEEE develops its standards through a consensus development process, approved by the American National Standards Institute, which brings together volunteers representing varied viewpoints and interests to achieve the final product. Volunteers are not necessarily members of the Institute and serve without compensation. While the IEEE administers the process and establishes rules to promote fairness in the consensus development process, the IEEE does not independently evaluate, test, or verify the accuracy of any of the information contained in its standards.

Use of an IEEE Standard is wholly voluntary. The IEEE disclaims liability for any personal injury, property, or other damage, of any nature whatsoever, whether special, indirect, consequential, or compensatory, directly or indirectly resulting from the publication, use of, or reliance upon this, or any other IEEE Standard document.

The IEEE does not warrant or represent the accuracy or content of the material contained herein, and expressly disclaims any express or implied warranty, including any implied warranty of merchantability or fitness for a specific purpose, or that the use of the material contained herein is free from patent infringement. IEEE Standards documents are supplied "AS IS".

The existence of an IEEE Standard does not imply that there are no other ways to produce, test, measure, purchase, market, or provide other goods and services related to the scope of the IEEE Standard. Furthermore, the viewpoint expressed at the time a standard is approved and issued is subject to change brought about through developments in the state of the art and comments received from users of the standard. Every IEEE Standard is subjected to review at least every five years for revision or reaffirmation. When a document is more than five years old and has not been reaffirmed, it is reasonable to conclude that its contents, although still of some value, do not wholly reflect the present state of the art. Users are cautioned to check to determine that they have the latest edition of any IEEE Standard.

In publishing and making this document available, the IEEE is not suggesting or rendering professional or other services for, or on behalf of, any person or entity. Nor is the IEEE undertaking to perform any duty owed by any other person or entity to another. Any person utilizing this, and any other IEEE Standards document, should rely upon the advice of a competent professional in determining the exercise of reasonable care in any given circumstances.

Interpretations: Occasionally questions may arise regarding the meaning of portions of standards as they relate to specific applications. When the need for interpretations is brought to the attention of the IEEE, the Institute will initiate action to prepare appropriate responses. Since IEEE Standards represent a consensus of concerned interests, it is important to ensure that any interpretation has also received the concurrence of a balance of interests. For this reason, IEEE and the members of its societies and Standards Coordinating Committees are not able to provide an instant response to interpretation requests except in those cases where the matter has previously received formal consideration.

Comments for revision of IEEE Standards are welcome from any interested party, regardless of membership affiliation with the IEEE.¹ Suggestions for changes in documents should be in the form of a proposed change of text, together with appropriate supporting comments. Comments on standards and requests for interpretations should be addressed to:

Secretary, IEEE-SA Standards Board, 445 Hoes Lane, P.O. Box 1331, Piscataway, NJ 08855-1331, U.S.A.

Attention is called to the possibility that implementation of this standard may require use of subject matter covered by patent rights. By publication of this standard, no position is taken with respect to the existence or validity of any patent rights in connection therewith. The IEEE shall not be responsible for identifying patents for which a license may be required by an IEEE Standard or for conducting inquiries into the legal validity or scope of those patents that are brought to its attention.

A patent holder has filed a statement of assurance that it will grant licenses under these rights without compensation or under reasonable rates and non-discriminatory, reasonable terms and conditions to all applicants desiring to obtain such licenses. The IEEE makes no representation as to the reasonableness of rates and/or terms and conditions of the license agreements offered by patent holders. Further information may be obtained from the IEEE Standards Department.

The IEEE and its designees are the sole entities that may authorize the use of IEEE-owned certification marks and/or trademarks to indicate compliance with the materials set forth herein. Authorization to photocopy portions of any individual standard for internal or personal use is granted in the U.S. by the Institute of Electrical and Electronics Engineers, Inc., provided that the appropriate fee is paid to the Copyright Clearance Center.² Permission to photocopy portions of any individual standard for educational classroom use can also be obtained through the Copyright Clearance Center. To arrange for payment of the licensing fee, please contact:

Copyright Clearance Center, Customer Service, 222 Rosewood Drive, Danvers, MA 01923, U.S.A., Tel.: +1 978 750 8400

Amendments, corrigenda, and interpretations for this standard, or information about the IEEE standards development process, may be found at <http://standards.ieee.org>.

Full catalog and ordering information on all IEEE publications is available from the IEEE Online Catalog & Store at <http://shop.ieee.org/store>.

1. For this standard, please send comments via the Austin Group as requested on page iii.

2. Please refer to the special provisions for this standard on page iii concerning permissions from both copyright holders and arrangements to cover photocopying and reproduction across the world, as well as by commercial organizations wishing to license the material for use in product documentation.

The Open Group

The Open Group, a vendor and technology-neutral consortium, is committed to delivering greater business efficiency by bringing together buyers and suppliers of information technology to lower the time, cost, and risks associated with integrating new technology across the enterprise.

The Open Group's mission is to offer all organizations concerned with open information infrastructures a forum to share knowledge, integrate open initiatives, and certify approved products and processes in a manner in which they continue to trust our impartiality.

In the global eCommerce world of today, no single economic entity can achieve independence while still ensuring interoperability. The assurance that products will interoperate with each other across differing systems and platforms is essential to the success of eCommerce and business workflow. The Open Group, with its proven testing and certification program, is the international guarantor of interoperability in the new century.

The Open Group provides opportunities to exchange information and shape the future of IT. The Open Group's members include some of the largest and most influential organizations in the world. The flexible structure of The Open Groups membership allows for almost any organization, no matter what their size, to join and have a voice in shaping the future of the IT world.

More information is available on The Open Group web site at <http://www.opengroup.org>.

The Open Group has over 15 years' experience in developing and operating certification programs and has extensive experience developing and facilitating industry adoption of test suites used to validate conformance to an open standard or specification. The Open Group portfolio of test suites includes the *Westwood* family of tests for this standard and the associated certification program for Version 3 of the Single UNIX Specification, as well tests for CDE, CORBA, Motif, Linux, LDAP, POSIX.1, POSIX.2, POSIX Realtime, Sockets, UNIX, XPG4, XNFS, XTI, and X11. The Open Group test tools are essential for proper development and maintenance of standards-based products, ensuring conformance of products to industry-standard APIs, applications portability, and interoperability. In-depth testing identifies defects at the earliest possible point in the development cycle, saving costs in development and quality assurance.

More information is available at <http://www.opengroup.org/testing>.

The Open Group publishes a wide range of technical documentation, the main part of which is focused on development of Technical and Product Standards and Guides, but which also includes white papers, technical studies, branding and testing documentation, and business titles. Full details and a catalog are available at <http://www.opengroup.org/pubs>.

As with all *live* documents, Technical Standards and Specifications require revision to align with new developments and associated international standards. To distinguish between revised specifications which are fully backwards compatible and those which are not:

- A new *Version* indicates there is no change to the definitive information contained in the previous publication of that title, but additions/extensions are included. As such, it *replaces* the previous publication.
- A new *Issue* indicates there is substantive change to the definitive information contained in the previous publication of that title, and there may also be additions/extensions. As such, both previous and new documents are maintained as current publications.

Readers should note that Corrigenda may apply to any publication. Corrigenda information is published at <http://www.opengroup.org/corrigenda>.

Full catalog and ordering information on all Open Group publications is available at <http://www.opengroup.org/pubs>.

Contents

Part	A	Base Definitions.....	1
Appendix	A	Rationale for Base Definitions.....	3
	A.1	Introduction	3
	A.1.1	Scope.....	3
	A.1.2	Conformance	5
	A.1.3	Normative References	5
	A.1.4	Terminology	5
	A.1.5	Portability	8
	A.1.5.1	Codes	8
	A.1.5.2	Margin Code Notation.....	8
	A.2	Conformance	9
	A.2.1	Implementation Conformance.....	9
	A.2.1.1	Requirements.....	9
	A.2.1.2	Documentation.....	9
	A.2.1.3	POSIX Conformance	10
	A.2.1.4	XSI Conformance	10
	A.2.1.5	Option Groups.....	11
	A.2.1.6	Options.....	12
	A.2.2	Application Conformance.....	12
	A.2.2.1	Strictly Conforming POSIX Application.....	12
	A.2.2.2	Conforming POSIX Application.....	12
	A.2.2.3	Conforming POSIX Application Using Extensions.....	12
	A.2.2.4	Strictly Conforming XSI Application	12
	A.2.2.5	Conforming XSI Application Using Extensions.....	12
	A.2.3	Language-Dependent Services for the C Programming Language.....	13
	A.2.4	Other Language-Related Specifications.....	13
	A.3	Definitions	13
	A.4	General Concepts.....	33
	A.4.1	Concurrent Execution.....	33
	A.4.2	Directory Protection	33
	A.4.3	Extended Security Controls.....	33
	A.4.4	File Access Permissions	33
	A.4.5	File Hierarchy	34
	A.4.6	Filenames.....	34
	A.4.7	File Times Update.....	35
	A.4.8	Host and Network Byte Order.....	36
	A.4.9	Measurement of Execution Time.....	36
	A.4.10	Memory Synchronization.....	36
	A.4.11	Pathname Resolution	38
	A.4.12	Process ID Reuse	39
	A.4.13	Scheduling Policy.....	39

A.4.14	Seconds Since the Epoch	39
A.4.15	Semaphore.....	40
A.4.16	Thread-Safety.....	40
A.4.17	Tracing.....	40
A.4.18	Treatment of Error Conditions for Mathematical Functions	41
A.4.19	Treatment of NaN Arguments for Mathematical Functions	41
A.4.20	Utility.....	41
A.4.21	Variable Assignment	41
A.5	File Format Notation.....	41
A.6	Character Set.....	42
A.6.1	Portable Character Set.....	42
A.6.2	Character Encoding.....	42
A.6.3	C Language Wide-Character Codes	42
A.6.4	Character Set Description File.....	42
A.6.4.1	State-Dependent Character Encodings	42
A.7	Locale.....	45
A.7.1	General.....	45
A.7.2	POSIX Locale	45
A.7.3	Locale Definition.....	45
A.7.3.1	LC_CTYPE.....	46
A.7.3.2	LC_COLLATE.....	47
A.7.3.3	LC_MONETARY.....	49
A.7.3.4	LC_NUMERIC.....	50
A.7.3.5	LC_TIME.....	50
A.7.3.6	LC_MESSAGES	51
A.7.4	Locale Definition Grammar	52
A.7.4.1	Locale Lexical Conventions.....	52
A.7.4.2	Locale Grammar.....	52
A.7.5	Locale Definition Example.....	52
A.8	Environment Variables	55
A.8.1	Environment Variable Definition	55
A.8.2	Internationalization Variables.....	56
A.8.3	Other Environment Variables.....	57
A.9	Regular Expressions.....	58
A.9.1	Regular Expression Definitions	58
A.9.2	Regular Expression General Requirements.....	59
A.9.3	Basic Regular Expressions	60
A.9.3.1	BREs Matching a Single Character or Collating Element.....	60
A.9.3.2	BRE Ordinary Characters.....	60
A.9.3.3	BRE Special Characters.....	60
A.9.3.4	Periods in BREs.....	60
A.9.3.5	RE Bracket Expression	60
A.9.3.6	BREs Matching Multiple Characters.....	62
A.9.3.7	BRE Precedence	62
A.9.3.8	BRE Expression Anchoring.....	62
A.9.4	Extended Regular Expressions	63
A.9.4.1	EREs Matching a Single Character or Collating Element.....	63
A.9.4.2	ERE Ordinary Characters.....	63

A.9.4.3	ERE Special Characters.....	63
A.9.4.4	Periods in EREs.....	63
A.9.4.5	ERE Bracket Expression.....	64
A.9.4.6	EREs Matching Multiple Characters.....	64
A.9.4.7	ERE Alternation.....	64
A.9.4.8	ERE Precedence	64
A.9.4.9	ERE Expression Anchoring.....	64
A.9.5	Regular Expression Grammar.....	64
A.9.5.1	BRE/ERE Grammar Lexical Conventions.....	64
A.9.5.2	RE and Bracket Expression Grammar	65
A.9.5.3	ERE Grammar.....	65
A.10	Directory Structure and Devices	65
A.10.1	Directory Structure and Files	65
A.10.2	Output Devices and Terminal Types	65
A.11	General Terminal Interface	66
A.11.1	Interface Characteristics	67
A.11.1.1	Opening a Terminal Device File	67
A.11.1.2	Process Groups.....	67
A.11.1.3	The Controlling Terminal.....	68
A.11.1.4	Terminal Access Control	68
A.11.1.5	Input Processing and Reading Data.....	69
A.11.1.6	Canonical Mode Input Processing	69
A.11.1.7	Non-Canonical Mode Input Processing.....	69
A.11.1.8	Writing Data and Output Processing	70
A.11.1.9	Special Characters.....	70
A.11.1.10	Modem Disconnect.....	70
A.11.1.11	Closing a Terminal Device File	70
A.11.2	Parameters that Can be Set	70
A.11.2.1	The termios Structure	70
A.11.2.2	Input Modes.....	70
A.11.2.3	Output Modes	71
A.11.2.4	Control Modes.....	71
A.11.2.5	Local Modes	71
A.11.2.6	Special Control Characters	71
A.12	Utility Conventions.....	72
A.12.1	Utility Argument Syntax.....	72
A.12.2	Utility Syntax Guidelines	73
A.13	Headers	75
A.13.1	Format of Entries.....	75

Part	B	System Interfaces.....	77
Appendix	B	Rationale for System Interfaces.....	79
	B.1	Introduction.....	79
	B.1.1	Scope.....	79
	B.1.2	Conformance.....	79
	B.1.3	Normative References.....	79
	B.1.4	Change History.....	79
	B.1.5	Terminology.....	85
	B.1.6	Definitions.....	85
	B.1.7	Relationship to Other Formal Standards.....	85
	B.1.8	Portability.....	85
	B.1.8.1	Codes.....	85
	B.1.9	Format of Entries.....	85
	B.2	General Information.....	86
	B.2.1	Use and Implementation of Functions.....	86
	B.2.2	The Compilation Environment.....	87
	B.2.2.1	POSIX.1 Symbols.....	87
	B.2.2.2	The Name Space.....	88
	B.2.3	Error Numbers.....	91
	B.2.3.1	Additional Error Numbers.....	95
	B.2.4	Signal Concepts.....	95
	B.2.4.1	Signal Generation and Delivery.....	96
	B.2.4.2	Realtime Signal Generation and Delivery.....	98
	B.2.4.3	Signal Actions.....	101
	B.2.4.4	Signal Effects on Other Functions.....	104
	B.2.5	Standard I/O Streams.....	104
	B.2.5.1	Interaction of File Descriptors and Standard I/O Streams.....	104
	B.2.5.2	Stream Orientation and Encoding Rules.....	104
	B.2.6	STREAMS.....	104
	B.2.6.1	Accessing STREAMS.....	105
	B.2.7	XSI Interprocess Communication.....	105
	B.2.7.1	IPC General Information.....	105
	B.2.8	Realtime.....	106
	B.2.8.1	Realtime Signals.....	111
	B.2.8.2	Asynchronous I/O.....	113
	B.2.8.3	Memory Management.....	116
	B.2.8.4	Process Scheduling.....	129
	B.2.8.5	Clocks and Timers.....	135
	B.2.9	Threads.....	151
	B.2.9.1	Thread-Safety.....	164
	B.2.9.2	Thread IDs.....	167
	B.2.9.3	Thread Mutexes.....	168
	B.2.9.4	Thread Scheduling.....	168
	B.2.9.5	Thread Cancellation.....	172
	B.2.9.6	Thread Read-Write Locks.....	176
	B.2.9.7	Thread Interactions with Regular File Operations.....	178
	B.2.10	Sockets.....	178

B.2.10.1	Address Families.....	178
B.2.10.2	Addressing	178
B.2.10.3	Protocols	178
B.2.10.4	Routing.....	178
B.2.10.5	Interfaces.....	178
B.2.10.6	Socket Types.....	178
B.2.10.7	Socket I/O Mode.....	178
B.2.10.8	Socket Owner.....	179
B.2.10.9	Socket Queue Limits	179
B.2.10.10	Pending Error.....	179
B.2.10.11	Socket Receive Queue.....	179
B.2.10.12	Socket Out-of-Band Data State	179
B.2.10.13	Connection Indication Queue	179
B.2.10.14	Signals	179
B.2.10.15	Asynchronous Errors	179
B.2.10.16	Use of Options.....	179
B.2.10.17	Use of Sockets for Local UNIX Connections.....	179
B.2.10.18	Use of Sockets over Internet Protocols.....	179
B.2.10.19	Use of Sockets over Internet Protocols Based on IPv4.....	179
B.2.10.20	Use of Sockets over Internet Protocols Based on IPv6.....	179
B.2.11	Tracing.....	180
B.2.11.1	Objectives	180
B.2.11.2	Trace Model.....	185
B.2.11.3	Trace Programming Examples	190
B.2.11.4	Rationale on Trace for Debugging.....	198
B.2.11.5	Rationale on Trace Event Type Name Space.....	198
B.2.11.6	Rationale on Trace Events Type Filtering	200
B.2.11.7	Tracing, pthread API.....	202
B.2.11.8	Rationale on Triggering.....	203
B.2.11.9	Rationale on Timestamp Clock.....	203
B.2.11.10	Rationale on Different Overrun Conditions.....	204
B.2.12	Data Types.....	204
B.3	System Interfaces	207
B.3.1	Examples for Spawn.....	207
Part	C Shell and Utilities.....	217
Appendix	C Rationale for Shell and Utilities	219
C.1	Introduction	219
C.1.1	Scope.....	219
C.1.2	Conformance	219
C.1.3	Normative References	219
C.1.4	Change History	219
C.1.5	Terminology	220
C.1.6	Definitions.....	220
C.1.7	Relationship to Other Documents.....	220
C.1.7.1	System Interfaces	220
C.1.7.2	Concepts Derived from the ISO C Standard.....	221

C.1.8	Portability	221
C.1.8.1	Codes	221
C.1.9	Utility Limits	222
C.1.10	Grammar Conventions	224
C.1.11	Utility Description Defaults	225
C.1.12	Considerations for Utilities in Support of Files of Arbitrary Size ..	228
C.1.13	Built-In Utilities	229
C.2	Shell Command Language	230
C.2.1	Shell Introduction	230
C.2.2	Quoting	231
C.2.2.1	Escape Character (Backslash)	231
C.2.2.2	Single-Quotes	231
C.2.2.3	Double-Quotes	231
C.2.3	Token Recognition	232
C.2.3.1	Alias Substitution	233
C.2.4	Reserved Words	233
C.2.5	Parameters and Variables	234
C.2.5.1	Positional Parameters	234
C.2.5.2	Special Parameters	234
C.2.5.3	Shell Variables	235
C.2.6	Word Expansions	236
C.2.6.1	Tilde Expansion	237
C.2.6.2	Parameter Expansion	238
C.2.6.3	Command Substitution	238
C.2.6.4	Arithmetic Expansion	239
C.2.6.5	Field Splitting	241
C.2.6.6	Pathname Expansion	241
C.2.6.7	Quote Removal	241
C.2.7	Redirection	241
C.2.7.1	Redirecting Input	243
C.2.7.2	Redirecting Output	243
C.2.7.3	Appending Redirected Output	243
C.2.7.4	Here-Document	243
C.2.7.5	Duplicating an Input File Descriptor	243
C.2.7.6	Duplicating an Output File Descriptor	243
C.2.7.7	Open File Descriptors for Reading and Writing	243
C.2.8	Exit Status and Errors	243
C.2.8.1	Consequences of Shell Errors	243
C.2.8.2	Exit Status for Commands	243
C.2.9	Shell Commands	244
C.2.9.1	Simple Commands	244
C.2.9.2	Pipelines	246
C.2.9.3	Lists	247
C.2.9.4	Compound Commands	248
C.2.9.5	Function Definition Command	249
C.2.10	Shell Grammar	251
C.2.10.1	Shell Grammar Lexical Conventions	252
C.2.10.2	Shell Grammar Rules	252

C.2.11	Signals and Error Handling	252	
C.2.12	Shell Execution Environment	252	
C.2.13	Pattern Matching Notation	252	
C.2.13.1	Patterns Matching a Single Character	252	
C.2.13.2	Patterns Matching Multiple Characters	253	
C.2.13.3	Patterns Used for Filename Expansion	253	
C.2.14	Special Built-In Utilities	254	
C.3	Batch Environment Services and Utilities	254	
C.3.1	Batch General Concepts	257	
C.3.2	Batch Services	259	
C.3.3	Common Behavior for Batch Environment Utilities	260	
C.4	Utilities	260	
Part	D	Portability Considerations	265
Appendix	D	Portability Considerations (Informative)	267
D.1	User Requirements	267	
D.1.1	Configuration Interrogation	268	
D.1.2	Process Management	268	
D.1.3	Access to Data	268	
D.1.4	Access to the Environment	268	
D.1.5	Access to Determinism and Performance Enhancements	268	
D.1.6	Operating System-Dependent Profile	268	
D.1.7	I/O Interaction	268	
D.1.8	Internationalization Interaction	269	
D.1.9	C-Language Extensions	269	
D.1.10	Command Language	269	
D.1.11	Interactive Facilities	269	
D.1.12	Accomplish Multiple Tasks Simultaneously	269	
D.1.13	Complex Data Manipulation	269	
D.1.14	File Hierarchy Manipulation	269	
D.1.15	Locale Configuration	269	
D.1.16	Inter-User Communication	270	
D.1.17	System Environment	270	
D.1.18	Printing	270	
D.1.19	Software Development	270	
D.2	Portability Capabilities	270	
D.2.1	Configuration Interrogation	271	
D.2.2	Process Management	271	
D.2.3	Access to Data	272	
D.2.4	Access to the Environment	272	
D.2.5	Bounded (Realtime) Response	273	
D.2.6	Operating System-Dependent Profile	273	
D.2.7	I/O Interaction	273	
D.2.8	Internationalization Interaction	273	
D.2.9	C-Language Extensions	274	
D.2.10	Command Language	274	
D.2.11	Interactive Facilities	274	

D.2.12	Accomplish Multiple Tasks Simultaneously	275
D.2.13	Complex Data Manipulation.....	275
D.2.14	File Hierarchy Manipulation	275
D.2.15	Locale Configuration.....	276
D.2.16	Inter-User Communication.....	276
D.2.17	System Environment.....	276
D.2.18	Printing.....	276
D.2.19	Software Development	277
D.2.20	Future Growth.....	277
D.3	Profiling Considerations.....	277
D.3.1	Configuration Options.....	277
D.3.2	Configuration Options (Shell and Utilities)	278
D.3.3	Configurable Limits.....	279
D.3.4	Configuration Options (System Interfaces)	280
D.3.5	Configurable Limits.....	285
D.3.6	Optional Behavior.....	288
Part E	Subprofiling Considerations	289
Appendix E	Subprofiling Considerations (Informative).....	291
E.1	Subprofiling Option Groups.....	291
	Index.....	297
List of Figures		
B-1	Example of a System with Typed Memory	124
B-2	Trace System Overview: for Offline Analysis.....	185
B-3	Trace System Overview: for Online Analysis	186
B-4	Trace System Overview: States of a Trace Stream	188
B-5	Trace Another Process.....	198
B-6	Trace Name Space Overview: With Third-Party Library	199
List of Tables		
A-1	Historical Practice for Symbolic Links	30

Introduction

Note: This introduction is not part of IEEE Std 1003.1-2001, Standard for Information Technology — Portable Operating System Interface (POSIX).

This standard has been jointly developed by the IEEE and The Open Group. It is both an IEEE Standard and an Open Group Technical Standard.

The Austin Group

This standard was developed, and is maintained, by a joint working group of members of the IEEE Portable Applications Standards Committee, members of The Open Group, and members of ISO/IEC Joint Technical Committee 1. This joint working group is known as the Austin Group.³ The Austin Group arose out of discussions amongst the parties which started in early 1998, leading to an initial meeting and formation of the group in September 1998. The purpose of the Austin Group has been to revise, combine, and update the following standards: ISO/IEC 9945-1, ISO/IEC 9945-2, IEEE Std 1003.1, IEEE Std 1003.2, and the Base Specifications of The Open Group Single UNIX Specification.

After two initial meetings, an agreement was signed in July 1999 between The Open Group and the Institute of Electrical and Electronics Engineers (IEEE), Inc., to formalize the project with the first draft of the revised specifications being made available at the same time. Under this agreement, The Open Group and IEEE agreed to share joint copyright of the resulting work. The Open Group has provided the chair and secretariat for the Austin Group.

The base document for the revision was The Open Group's Base volumes of its Single UNIX Specification, Version 2. These were selected since they were a superset of the existing POSIX.1 and POSIX.2 specifications and had some organizational aspects that would benefit the audience for the new revision.

The approach to specification development has been one of “write once, adopt everywhere”, with the deliverables being a set of specifications that carry the IEEE POSIX designation and The Open Group's Technical Standard designation, and, if approved, an ISO/IEC designation. This set of specifications forms the core of the Single UNIX Specification, Version 3.

This unique development has combined both the industry-led efforts and the formal standardization activities into a single initiative, and included a wide spectrum of participants. The Austin Group continues as the maintenance body for this document.

Anyone wishing to participate in the Austin Group should contact the chair with their request. There are no fees for participation or membership. You may participate as an observer or as a contributor. You do not have to attend face-to-face meetings to participate; electronic participation is most welcome. For more information on the Austin Group and how to participate, see <http://www.opengroup.org/austin>.

3. The Austin Group is named after the location of the inaugural meeting held at the IBM facility in Austin, Texas in September 1998.

Background

The developers of this standard represent a cross section of hardware manufacturers, vendors of operating systems and other software development tools, software designers, consultants, academics, authors, applications programmers, and others.

Conceptually, this standard describes a set of fundamental services needed for the efficient construction of application programs. Access to these services has been provided by defining an interface, using the C programming language, a command interpreter, and common utility programs that establish standard semantics and syntax. Since this interface enables application writers to write portable applications—it was developed with that goal in mind—it has been designated POSIX,⁴ an acronym for Portable Operating System Interface.

Although originated to refer to the original IEEE Std 1003.1-1988, the name POSIX more correctly refers to a *family* of related standards: IEEE Std 1003.*n* and the parts of ISO/IEC 9945. In earlier editions of the IEEE standard, the term POSIX was used as a synonym for IEEE Std 1003.1-1988. A preferred term, POSIX.1, emerged. This maintained the advantages of readability of the symbol “POSIX” without being ambiguous with the POSIX family of standards.

Audience

The intended audience for this standard is all persons concerned with an industry-wide standard operating system based on the UNIX system. This includes at least four groups of people:

1. Persons buying hardware and software systems
2. Persons managing companies that are deciding on future corporate computing directions
3. Persons implementing operating systems, and especially
4. Persons developing applications where portability is an objective

Purpose

Several principles guided the development of this standard:

- Application-Oriented

The basic goal was to promote portability of application programs across UNIX system environments by developing a clear, consistent, and unambiguous standard for the interface specification of a portable operating system based on the UNIX system documentation. This standard codifies the common, existing definition of the UNIX system.

- Interface, Not Implementation

This standard defines an interface, not an implementation. No distinction is made between library functions and system calls; both are referred to as functions. No details of the implementation of any function are given (although historical practice is sometimes indicated in the RATIONALE section). Symbolic names are given for constants (such as signals and error numbers) rather than numbers.

4. The name POSIX was suggested by Richard Stallman. It is expected to be pronounced *pahz-icks*, as in *positive*, not *poh-six*, or other variations. The pronunciation has been published in an attempt to promulgate a standardized way of referring to a standard operating system interface.

- Source, Not Object, Portability

This standard has been written so that a program written and translated for execution on one conforming implementation may also be translated for execution on another conforming implementation. This standard does not guarantee that executable (object or binary) code will execute under a different conforming implementation than that for which it was translated, even if the underlying hardware is identical.

- The C Language

The system interfaces and header definitions are written in terms of the standard C language as specified in the ISO C standard.

- No Superuser, No System Administration

There was no intention to specify all aspects of an operating system. System administration facilities and functions are excluded from this standard, and functions usable only by the superuser have not been included. Still, an implementation of the standard interface may also implement features not in this standard. This standard is also not concerned with hardware constraints or system maintenance.

- Minimal Interface, Minimally Defined

In keeping with the historical design principles of the UNIX system, the mandatory core facilities of this standard have been kept as minimal as possible. Additional capabilities have been added as optional extensions.

- Broadly Implementable

The developers of this standard endeavored to make all specified functions implementable across a wide range of existing and potential systems, including:

1. All of the current major systems that are ultimately derived from the original UNIX system code (Version 7 or later)
2. Compatible systems that are not derived from the original UNIX system code
3. Emulations hosted on entirely different operating systems
4. Networked systems
5. Distributed systems
6. Systems running on a broad range of hardware

No direct references to this goal appear in this standard, but some results of it are mentioned in the Rationale (Informative) volume.

- Minimal Changes to Historical Implementations

When the original version of IEEE Std 1003.1 was published, there were no known historical implementations that did not have to change. However, there was a broad consensus on a set of functions, types, definitions, and concepts that formed an interface that was common to most historical implementations.

The adoption of the 1988 and 1990 IEEE system interface standards, the 1992 IEEE shell and utilities standard, the various Open Group (formerly X/Open) specifications, and the subsequent revisions and addenda to all of them have consolidated this consensus, and this revision reflects the significantly increased level of consensus arrived at since the original versions. The earlier standards and their modifications specified a number of areas where consensus had not been reached before, and these are now reflected in this revision. The authors of the original versions tried, as much as possible, to follow the principles below

when creating new specifications:

1. By standardizing an interface like one in an historical implementation; for example, directories
2. By specifying an interface that is readily implementable in terms of, and backwards-compatible with, historical implementations, such as the extended *tar* format defined in the *pax* utility
3. By specifying an interface that, when added to an historical implementation, will not conflict with it; for example, the *sigaction()* function

This revision tries to minimize the number of changes required to implementations which conform to the earlier versions of the approved standards to bring them into conformance with the current standard. Specifically, the scope of this work excluded doing any “new” work, but rather collecting into a single document what had been spread across a number of documents, and presenting it in what had been proven in practice to be a more effective way. Some changes to prior conforming implementations were unavoidable, primarily as a consequence of resolving conflicts found in prior revisions, or which became apparent when bringing the various pieces together.

However, since it references the 1999 version of the ISO C standard, and no longer supports “Common Usage C”, there are a number of unavoidable changes. Applications portability is similarly affected.

This standard is specifically not a codification of a particular vendor’s product.

It should be noted that implementations will have different kinds of extensions. Some will reflect “historical usage” and will be preserved for execution of pre-existing applications. These functions should be considered “obsolescent” and the standard functions used for new applications. Some extensions will represent functions beyond the scope of this standard. These need to be used with careful management to be able to adapt to future extensions of this standard and/or port to implementations that provide these services in a different manner.

- Minimal Changes to Existing Application Code

A goal of this standard was to minimize additional work for the developers of applications. However, because every known historical implementation will have to change at least slightly to conform, some applications will have to change.

This Standard

This standard defines the Portable Operating System Interface (POSIX) requirements and consists of the following volumes:

- Base Definitions
- Shell and Utilities
- System Interfaces
- Rationale (Informative) (this volume)

This Volume

This volume is being published to assist in the process of review. It contains historical information concerning the contents of this standard and why features were included or discarded by the standard developers. It also contains notes of interest to application programmers on recommended programming practices, emphasizing the consequences of some aspects of this standard that may not be immediately apparent.

This volume is organized in parallel to the normative volumes of this standard, with a separate part for each of the three normative volumes.

Within this volume, the following terms are used:

base standard

The portions of this standard that are not optional, equivalent to the definitions of *classic* POSIX.1 and POSIX.2.

POSIX.0

Although this term is not used in the normative text of this standard, it is used in this volume to refer to IEEE Std 1003.0-1995.

POSIX.1b

Although this term is not used in the normative text of this standard, it is used in this volume to refer to the elements of the POSIX Realtime Extension amendment. (This was earlier referred to as POSIX.4 during the standard development process.)

POSIX.1c

Although this term is not used in the normative text of this standard, it is used in this volume to refer to the POSIX Threads Extension amendment. (This was earlier referred to as POSIX.4a during the standard development process.)

standard developers

The individuals and companies in the development organizations responsible for this standard: the IEEE P1003.1 working groups, The Open Group Base working group, advised by the hundreds of individual technical experts who balloted the draft standards within the Austin Group, and the member bodies and technical experts of ISO/IEC JTC 1/SC22/WG15.

XSI extension

The portions of this standard addressing the extension added for support of the Single UNIX Specification.

Participants

The Austin Group

This standard was prepared by the Austin Group, sponsored by the Portable Applications Standards Committee of the IEEE Computer Society, The Open Group, and ISO/SC22 WG15. At the time of approval, the membership of the Austin Group was as follows.

Andrew Josey, Chair

Donald W. Cragun, Organizational Representative, IEEE PASC

Nicholas Stoughton, Organizational Representative, ISO/SC22 WG15

Mark Brown, Organizational Representative, The Open Group

Cathy Hughes, Technical Editor

Austin Group Technical Reviewers

Peter Anvin

Bouazza Bachar

Theodore P. Baker

Walter Briscoe

Mark Brown

Dave Butenhof

Geoff Clare

Donald W. Cragun

Lee Damico

Ulrich Drepper

Paul Eggert

Joanna Farley

Clive D.W. Feather

Andrew Gollan

Michael Gonzalez

Joseph M. Gwinn

Jon Hitchcock

Yvette Ho Sang

Cathy Hughes

Lowell G. Johnson

Andrew Josey

Michael Kavanaugh

David Korn

Marc Aurele La France

Jim Meyering

Gary Miller

Finnbarr P. Murphy

Joseph S. Myers

Sandra O'Donnell

Frank Prindle

Curtis Royster Jr.

Glen Seeds

Keld Jorn Simonsen

Raja Srinivasan

Nicholas Stoughton

Donn S. Terry

Fred Tydeman

Peter Van Der Veen

James Youngman

Jim Zepeda

Jason Zions

Participants

Austin Group Working Group Members

Harold C. Adams
Peter Anvin
Pierre-Jean Arcos
Jay Ashford
Bouazza Bachar
Theodore P. Baker
Robert Barned
Joel Berman
David J. Blackwood
Shirley Bockstahler-Brandt
James Bottomley
Walter Briscoe
Andries Brouwer
Mark Brown
Eric W. Burger
Alan Burns
Andries Brouwer
Dave Butenhof
Keith Chow
Geoff Clare
Donald W. Cragun
Lee Damico
Juan Antonio De La Puente
Ming De Zhou
Steven J. Dovich
Richard P. Draves
Ulrich Drepper
Paul Eggert
Philip H. Enslow
Joanna Farley
Clive D.W. Feather
Pete Forman
Mark Funkenhauser
Lois Goldthwaite
Andrew Gollan

Michael Gonzalez
Karen D. Gordon
Joseph M. Gwinn
Steven A. Haaser
Charles E. Hammons
Chris J. Harding
Barry Hedquist
Vincent E. Henley
Karl Heubaum
Jon Hitchcock
Yvette Ho Sang
Niklas Holsti
Thomas Hosmer
Cathy Hughes
Jim D. Isaak
Lowell G. Johnson
Michael B. Jones
Andrew Josey
Michael J. Karels
Michael Kavanaugh
David Korn
Steven Kramer
Thomas M. Kurihara
Marc Aurele La France
C. Douglass Locke
Nick Maclaren
Roger J. Martin
Craig H. Meyer
Jim Meyering
Gary Miller
Finnbarr P. Murphy
Joseph S. Myers
John Napier
Peter E. Obermayer
James T. Oblinger

Sandra O'Donnell
Frank Prindle
Francois Riche
John D. Riley
Andrew K. Roach
Helmut Roth
Jaideep Roy
Curtis Royster Jr.
Stephen C. Schwarm
Glen Seeds
Richard Seibel
David L. Shroads Jr.
W. Olin Sibert
Keld Jorn Simonsen
Curtis Smith
Raja Srinivasan
Nicholas Stoughton
Marc J. Teller
Donn S. Terry
Fred Tydeman
Mark-Rene Uchida
Scott A. Valcourt
Peter Van Der Veen
Michael W. Vannier
Eric Vought
Frederick N. Webb
Paul A.T. Wolfgang
Garrett Wollman
James Youngman
Oren Yuen
Janusz Zalewski
Jim Zepeda
Jason Zions

The Open Group

When The Open Group approved the Base Specifications, Issue 6 on 12 September 2001, the membership of The Open Group Base Working Group was as follows.

Andrew Josey, Chair

Finnbarr P. Murphy, Vice-Chair

Mark Brown, Austin Group Liaison

Cathy Hughes, Technical Editor

Base Working Group Members

Bouazza Bachar

Mark Brown

Dave Butenhof

Donald W. Cragun

Larry Dwyer

Joanna Farley

Andrew Gollan

Karen D. Gordon

Gary Miller

Finnbarr P. Murphy

Frank Prindle

Andrew K. Roach

Curtis Royster Jr.

Nicholas Stoughton

Kenjiro Tsuji

Participants

IEEE

When the IEEE Standards Board approved IEEE Std 1003.1-2001 on 6 December 2001, the membership of the committees was as follows.

Portable Applications Standards Committee (PASC)

Lowell G. Johnson, Chair

Joseph M. Gwinn, Vice-Chair

Jay Ashford, Functional Chair

Andrew Josey, Functional Chair

Curtis Royster Jr., Functional Chair

Nicholas Stoughton, Secretary

Balloting Committee

The following members of the balloting committee voted on IEEE Std 1003.1-2001. Balloters may have voted for approval, disapproval, or abstention.

Harold C. Adams

Pierre-Jean Arcos

Jay Ashford

Theodore P. Baker

Robert Barned

David J. Blackwood

Shirley Bockstahler-Brandt

James Bottomley

Mark Brown

Eric W. Burger

Alan Burns

Dave Butenhof

Keith Chow

Donald W. Cragun

Juan Antonio De La Puente

Ming De Zhou

Steven J. Dovich

Richard P. Draves

Philip H. Enslow

Michael Gonzalez

Karen D. Gordon

Joseph M. Gwinn

Steven A. Haaser

Charles E. Hammons

Chris J. Harding

Barry Hedquist

Vincent E. Henley

Karl Heubaum

Niklas Holsti

Thomas Hosmer

Jim D. Isaak

Lowell G. Johnson

Michael B. Jones

Andrew Josey

Michael J. Karels

Steven Kramer

Thomas M. Kurihara

C. Douglass Locke

Roger J. Martin

Craig H. Meyer

Finnbarr P. Murphy

John Napier

Peter E. Obermayer

James T. Oblinger

Frank Prindle

Francois Riche

John D. Riley

Andrew K. Roach

Helmut Roth

Jaideep Roy

Curtis Royster Jr.

Stephen C. Schwarm

Richard Seibel

David L. Shroads Jr.

W. Olin Sibert

Keld Jorn Simonsen

Nicholas Stoughton

Donn S. Terry

Mark-Rene Uchida

Scott A. Valcourt

Michael W. Vannier

Frederick N. Webb

Paul A.T. Wolfgang

Oren Yuen

Janusz Zalewski

The following organizational representative voted on this standard:

Andrew Josey, X/Open Company Ltd.

IEEE-SA Standards Board

When the IEEE-SA Standards Board approved IEEE Std 1003.1-2001 on 6 December 2001, it had the following membership:

Donald N. Heirman, Chair

James T. Carlo, Vice-Chair

Judith Gorman, Secretary

Satish K. Aggarwal

Mark D. Bowman

Gary R. Engmann

Harold E. Epstein

H. Landis Floyd

Jay Forster*

Howard M. Frazier

Ruben D. Garzon

James H. Gurney

Richard J. Holleman

Lowell G. Johnson

Robert J. Kennelly

Joseph L. Koepfinger*

Peter H. Lips

L. Bruce McClung

Daleep C. Mohla

James W. Moore

Robert F. Munzner

Ronald C. Petersen

Gerald H. Peterson

John B. Posey

Gary S. Robinson

Akio Tojo

Donald W. Zipse

Also included are the following non-voting IEEE-SA Standards Board liaisons:

Alan Cookson, NIST Representative

Donald R. Volzka, TAB Representative

Yvette Ho Sang, **Don Messina**, **Savoula Amanatidis**, IEEE Project Editors

* Member Emeritus

Trademarks

The following information is given for the convenience of users of this standard and does not constitute endorsement of these products by The Open Group or the IEEE. There may be other products mentioned in the text that might be covered by trademark protection and readers are advised to verify them independently.

1003.1™ is a trademark of the Institute of Electrical and Electronic Engineers, Inc.

AIX® is a registered trademark of IBM Corporation.

AT&T® is a registered trademark of AT&T in the U.S.A. and other countries.

BSD™ is a trademark of the University of California, Berkeley, U.S.A.

Hewlett-Packard®, HP®, and HP-UX® are registered trademarks of Hewlett-Packard Company.

IBM® is a registered trademark of International Business Machines Corporation.

Motif®, OSF/1®, UNIX®, and the “X Device” are registered trademarks and IT DialTone™ and The Open Group™ are trademarks of The Open Group in the U.S. and other countries.

POSIX® is a registered trademark of the Institute of Electrical and Electronic Engineers, Inc.

Sun® and Sun Microsystems® are registered trademarks of Sun Microsystems, Inc.

/usr/group® is a registered trademark of UniForum, the International Network of UNIX System Users.

Acknowledgements

The contributions of the following organizations to the development of IEEE Std 1003.1-2001 are gratefully acknowledged:

- AT&T for permission to reproduce portions of its copyrighted System V Interface Definition (SVID) and material from the UNIX System V Release 2.0 documentation.
- The SC22 WG14 Committees.

This standard was prepared by the Austin Group, a joint working group of the IEEE, The Open Group, and ISO SC22 WG15.

Referenced Documents

Normative References

Normative references for this standard are defined in the Base Definitions volume.

Informative References

The following documents are referenced in this standard:

1984 /usr/group Standard

/usr/group Standards Committee, Santa Clara, CA, UniForum 1984.

Almasi and Gottlieb

George S. Almasi and Allan Gottlieb, *Highly Parallel Computing*, The Benjamin/Cummings Publishing Company, Inc., 1989, ISBN: 0-8053-0177-1.

ANSI C

American National Standard for Information Systems: Standard X3.159-1989, Programming Language C.

ANSI X3.226-1994

American National Standard for Information Systems: Standard X3.226-1994, Programming Language Common LISP.

Brawer

Steven Brawer, *Introduction to Parallel Programming*, Academic Press, 1989, ISBN: 0-12-128470-0.

DeRemer and Pennello Article

DeRemer, Frank and Pennello, Thomas J., *Efficient Computation of LALR(1) Look-Ahead Sets*, SigPlan Notices, Volume 15, No. 8, August 1979.

Draft ANSI X3J11.1

IEEE Floating Point draft report of ANSI X3J11.1 (NCEG).

FIPS 151-1

Federal Information Procurement Standard (FIPS) 151-1. Portable Operating System Interface (POSIX)—Part 1: System Application Program Interface (API) [C Language].

FIPS 151-2

Federal Information Procurement Standards (FIPS) 151-2, Portable Operating System Interface (POSIX)—Part 1: System Application Program Interface (API) [C Language].

HP-UX Manual

Hewlett-Packard HP-UX Release 9.0 Reference Manual, Third Edition, August 1992.

IEC 60559: 1989

IEC 60559: 1989, Binary Floating-Point Arithmetic for Microprocessor Systems (previously designated IEC 559: 1989).

IEEE Std 754-1985

IEEE Std 754-1985, IEEE Standard for Binary Floating-Point Arithmetic.

IEEE Std 854-1987

IEEE Std 854-1987, IEEE Standard for Radix-Independent Floating-Point Arithmetic.

- IEEE Std 1003.9-1992
IEEE Std 1003.9-1992, IEEE Standard for Information Technology — POSIX FORTRAN 77 Language Interfaces — Part 1: Binding for System Application Program Interface API.
- IETF RFC 791
Internet Protocol, Version 4 (IPv4), September 1981.
- IETF RFC 819
The Domain Naming Convention for Internet User Applications, Z. Su, J. Postel, August 1982.
- IETF RFC 822
Standard for the Format of ARPA Internet Text Messages, D.H. Crocker, August 1982.
- IETF RFC 919
Broadcasting Internet Datagrams, J. Mogul, October 1984.
- IETF RFC 920
Domain Requirements, J. Postel, J. Reynolds, October 1984.
- IETF RFC 921
Domain Name System Implementation Schedule, J. Postel, October 1984.
- IETF RFC 922
Broadcasting Internet Datagrams in the Presence of Subnets, J. Mogul, October 1984.
- IETF RFC 1034
Domain Names — Concepts and Facilities, P. Mockapetris, November 1987.
- IETF RFC 1035
Domain Names — Implementation and Specification, P. Mockapetris, November 1987.
- IETF RFC 1123
Requirements for Internet Hosts — Application and Support, R. Braden, October 1989.
- IETF RFC 1886
DNS Extensions to Support Internet Protocol, Version 6 (IPv6), C. Huitema, S. Thomson, December 1995.
- IETF RFC 2045
Multipurpose Internet Mail Extensions (MIME), Part 1: Format of Internet Message Bodies, N. Freed, N. Borenstein, November 1996.
- IETF RFC 2373
Internet Protocol, Version 6 (IPv6) Addressing Architecture, S. Deering, R. Hinden, July 1998.
- IETF RFC 2460
Internet Protocol, Version 6 (IPv6), S. Deering, R. Hinden, December 1998.
- Internationalisation Guide
Guide, July 1993, Internationalisation Guide, Version 2 (ISBN: 1-859120-02-4, G304), published by The Open Group.
- ISO C (1990)
ISO/IEC 9899:1990, Programming Languages — C, including Amendment 1:1995 (E), C Integrity (Multibyte Support Extensions (MSE) for ISO C).
- ISO 2375:1985
ISO 2375:1985, Data Processing — Procedure for Registration of Escape Sequences.

Referenced Documents

ISO 8652:1987

ISO 8652:1987, Programming Languages — Ada (technically identical to ANSI standard 1815A-1983).

ISO/IEC 1539:1990

ISO/IEC 1539:1990, Information Technology — Programming Languages — Fortran (technically identical to the ANSI X3.9-1978 standard [FORTRAN 77]).

ISO/IEC 4873:1991

ISO/IEC 4873:1991, Information Technology — ISO 8-bit Code for Information Interchange — Structure and Rules for Implementation.

ISO/IEC 6429:1992

ISO/IEC 6429:1992, Information Technology — Control Functions for Coded Character Sets.

ISO/IEC 6937:1994

ISO/IEC 6937:1994, Information Technology — Coded Character Set for Text Communication — Latin Alphabet.

ISO/IEC 8802-3:1996

ISO/IEC 8802-3:1996, Information Technology — Telecommunications and Information Exchange Between Systems — Local and Metropolitan Area Networks — Specific Requirements — Part 3: Carrier Sense Multiple Access with Collision Detection (CSMA/CD) Access Method and Physical Layer Specifications.

ISO/IEC 8859

ISO/IEC 8859, Information Technology — 8-Bit Single-Byte Coded Graphic Character Sets:

Part 1: Latin Alphabet No. 1

Part 2: Latin Alphabet No. 2

Part 3: Latin Alphabet No. 3

Part 4: Latin Alphabet No. 4

Part 5: Latin/Cyrillic Alphabet

Part 6: Latin/Arabic Alphabet

Part 7: Latin/Greek Alphabet

Part 8: Latin/Hebrew Alphabet

Part 9: Latin Alphabet No. 5

Part 10: Latin Alphabet No. 6

Part 13: Latin Alphabet No. 7

Part 14: Latin Alphabet No. 8

Part 15: Latin Alphabet No. 9

ISO POSIX-1:1996

ISO/IEC 9945-1:1996, Information Technology — Portable Operating System Interface (POSIX) — Part 1: System Application Program Interface (API) [C Language] (identical to ANSI/IEEE Std 1003.1-1996). Incorporating ANSI/IEEE Stds 1003.1-1990, 1003.1b-1993, 1003.1c-1995, and 1003.1i-1995.

ISO POSIX-2:1993

ISO/IEC 9945-2:1993, Information Technology — Portable Operating System Interface (POSIX) — Part 2: Shell and Utilities (identical to ANSI/IEEE Std 1003.2-1992, as amended by ANSI/IEEE Std 1003.2a-1992).

Issue 1

X/Open Portability Guide, July 1985 (ISBN: 0-444-87839-4).

Issue 2

X/Open Portability Guide, January 1987:

- Volume 1: XVS Commands and Utilities (ISBN: 0-444-70174-5)
- Volume 2: XVS System Calls and Libraries (ISBN: 0-444-70175-3)

Issue 3

X/Open Specification, 1988, 1989, February 1992:

- Commands and Utilities, Issue 3 (ISBN: 1-872630-36-7, C211); this specification was formerly X/Open Portability Guide, Issue 3, Volume 1, January 1989, XSI Commands and Utilities (ISBN: 0-13-685835-X, XO/XPG/89/002)
- System Interfaces and Headers, Issue 3 (ISBN: 1-872630-37-5, C212); this specification was formerly X/Open Portability Guide, Issue 3, Volume 2, January 1989, XSI System Interface and Headers (ISBN: 0-13-685843-0, XO/XPG/89/003)
- Curses Interface, Issue 3, contained in Supplementary Definitions, Issue 3 (ISBN: 1-872630-38-3, C213), Chapters 9 to 14 inclusive; this specification was formerly X/Open Portability Guide, Issue 3, Volume 3, January 1989, XSI Supplementary Definitions (ISBN: 0-13-685850-3, XO/XPG/89/004)
- Headers Interface, Issue 3, contained in Supplementary Definitions, Issue 3 (ISBN: 1-872630-38-3, C213), Chapter 19, Cpio and Tar Headers; this specification was formerly X/Open Portability Guide Issue 3, Volume 3, January 1989, XSI Supplementary Definitions (ISBN: 0-13-685850-3, XO/XPG/89/004)

Issue 4

CAE Specification, July 1992, published by The Open Group:

- System Interface Definitions (XBD), Issue 4 (ISBN: 1-872630-46-4, C204)
- Commands and Utilities (XCU), Issue 4 (ISBN: 1-872630-48-0, C203)
- System Interfaces and Headers (XSH), Issue 4 (ISBN: 1-872630-47-2, C202)

Issue 4, Version 2

CAE Specification, August 1994, published by The Open Group:

- System Interface Definitions (XBD), Issue 4, Version 2 (ISBN: 1-85912-036-9, C434)
- Commands and Utilities (XCU), Issue 4, Version 2 (ISBN: 1-85912-034-2, C436)
- System Interfaces and Headers (XSH), Issue 4, Version 2 (ISBN: 1-85912-037-7, C435)

Issue 5

Technical Standard, February 1997, published by The Open Group:

- System Interface Definitions (XBD), Issue 5 (ISBN: 1-85912-186-1, C605)
- Commands and Utilities (XCU), Issue 5 (ISBN: 1-85912-191-8, C604)
- System Interfaces and Headers (XSH), Issue 5 (ISBN: 1-85912-181-0, C606)

Knuth Article

Knuth, Donald E., *On the Translation of Languages from Left to Right*, Information and Control, Volume 8, No. 6, October 1965.

KornShell

Bolsky, Morris I. and Korn, David G., *The New KornShell Command and Programming Language*, March 1995, Prentice Hall.

Referenced Documents

MSE Working Draft

Working draft of ISO/IEC 9899:1990/Add3:Draft, Addendum 3 — Multibyte Support Extensions (MSE) as documented in the ISO Working Paper SC22/WG14/N205 dated 31 March 1992.

POSIX.0: 1995

IEEE Std 1003.0-1995, IEEE Guide to the POSIX Open System Environment (OSE) (identical to ISO/IEC TR 14252).

POSIX.1: 1988

IEEE Std 1003.1-1988, IEEE Standard for Information Technology — Portable Operating System Interface (POSIX) — Part 1: System Application Program Interface (API) [C Language].

POSIX.1: 1990

IEEE Std 1003.1-1990, IEEE Standard for Information Technology — Portable Operating System Interface (POSIX) — Part 1: System Application Program Interface (API) [C Language].

POSIX.1a

P1003.1a, Standard for Information Technology — Portable Operating System Interface (POSIX) — Part 1: System Application Program Interface (API) — (C Language) Amendment

POSIX.1d: 1999

IEEE Std 1003.1d-1999, IEEE Standard for Information Technology — Portable Operating System Interface (POSIX) — Part 1: System Application Program Interface (API) — Amendment 4: Additional Realtime Extensions [C Language].

POSIX.1g: 2000

IEEE Std 1003.1g-2000, IEEE Standard for Information Technology — Portable Operating System Interface (POSIX) — Part 1: System Application Program Interface (API) — Amendment 6: Protocol-Independent Interfaces (PII).

POSIX.1j: 2000

IEEE Std 1003.1j-2000, IEEE Standard for Information Technology — Portable Operating System Interface (POSIX) — Part 1: System Application Program Interface (API) — Amendment 5: Advanced Realtime Extensions [C Language].

POSIX.1q: 2000

IEEE Std 1003.1q-2000, IEEE Standard for Information Technology — Portable Operating System Interface (POSIX) — Part 1: System Application Program Interface (API) — Amendment 7: Tracing [C Language].

POSIX.2b

P1003.2b, Standard for Information Technology — Portable Operating System Interface (POSIX) — Part 2: Shell and Utilities — Amendment

POSIX.2d:-1994

IEEE Std 1003.2d: 1994, IEEE Standard for Information Technology — Portable Operating System Interface (POSIX) — Part 2: Shell and Utilities — Amendment 1: Batch Environment.

POSIX.13:-1998

IEEE Std 1003.13: 1998, IEEE Standard for Information Technology — Standardized Application Environment Profile (AEP) — POSIX Realtime Application Support.

Sarwate Article

Sarwate, Dilip V., *Computation of Cyclic Redundancy Checks via Table Lookup*, Communications of the ACM, Volume 30, No. 8, August 1988.

Sprunt, Sha, and Lehoczky

Sprunt, B., Sha, L., and Lehoczky, J.P., *Aperiodic Task Scheduling for Hard Real-Time Systems*, The Journal of Real-Time Systems, Volume 1, 1989, Pages 27-60.

SVID, Issue 1

American Telephone and Telegraph Company, System V Interface Definition (SVID), Issue 1; Morristown, NJ, UNIX Press, 1985.

SVID, Issue 2

American Telephone and Telegraph Company, System V Interface Definition (SVID), Issue 2; Morristown, NJ, UNIX Press, 1986.

SVID, Issue 3

American Telephone and Telegraph Company, System V Interface Definition (SVID), Issue 3; Morristown, NJ, UNIX Press, 1989.

The AWK Programming Language

Aho, Alfred V., Kernighan, Brian W., and Weinberger, Peter J., *The AWK Programming Language*, Reading, MA, Addison-Wesley 1988.

UNIX Programmer's Manual

American Telephone and Telegraph Company, *UNIX Time-Sharing System: UNIX Programmer's Manual*, 7th Edition, Murray Hill, NJ, Bell Telephone Laboratories, January 1979.

XNS, Issue 4

CAE Specification, August 1994, Networking Services, Issue 4 (ISBN: 1-85912-049-0, C438), published by The Open Group.

XNS, Issue 5

CAE Specification, February 1997, Networking Services, Issue 5 (ISBN: 1-85912-165-9, C523), published by The Open Group.

XNS, Issue 5.2

Technical Standard, January 2000, Networking Services (XNS), Issue 5.2 (ISBN: 1-85912-241-8, C808), published by The Open Group.

X/Open Curses, Issue 4, Version 2

CAE Specification, May 1996, X/Open Curses, Issue 4, Version 2 (ISBN: 1-85912-171-3, C610), published by The Open Group.

Yacc

Yacc: Yet Another Compiler Compiler, Stephen C. Johnson, 1978.

Source Documents

Parts of the following documents were used to create the base documents for this standard:

AIX 3.2 Manual

AIX Version 3.2 For RISC System/6000, Technical Reference: Base Operating System and Extensions, 1990, 1992 (Part No. SC23-2382-00).

OSF/1

OSF/1 Programmer's Reference, Release 1.2 (ISBN: 0-13-020579-6).

Referenced Documents

OSF AES

Application Environment Specification (AES) Operating System Programming Interfaces
Volume, Revision A (ISBN: 0-13-043522-8).

System V Release 2.0

- UNIX System V Release 2.0 Programmer's Reference Manual (April 1984 - Issue 2).
- UNIX System V Release 2.0 Programming Guide (April 1984 - Issue 2).

System V Release 4.2

Operating System API Reference, UNIX SVR4.2 (1992) (ISBN: 0-13-017658-3).

1 / *Rationale (Informative)*

2 **Part A:**

3 **Base Definitions**

4 *The Open Group*

5 *The Institute of Electrical and Electronics Engineers, Inc.*

Rationale for Base Definitions

6

7 A.1 Introduction

8 A.1.1 Scope

9 IEEE Std 1003.1-2001 is one of a family of standards known as POSIX. The family of standards
10 extends to many topics; IEEE Std 1003.1-2001 is known as POSIX.1 and consists of both
11 operating system interfaces and shell and utilities. IEEE Std 1003.1-2001 is technically identical
12 to The Open Group Base Specifications, Issue 6, which comprise the core volumes of the Single
13 UNIX Specification, Version 3.

14 Scope of IEEE Std 1003.1-2001

15 The (paraphrased) goals of this development were to produce a single common revision to the
16 overlapping POSIX.1 and POSIX.2 standards, and the Single UNIX Specification, Version 2. As
17 such, the scope of the revision includes the scopes of the original documents merged.

18 Since the revision includes merging the Base volumes of the Single UNIX Specification, many
19 features that were previously not “adopted” into earlier revisions of POSIX.1 and POSIX.2 are
20 now included in IEEE Std 1003.1-2001. In most cases, these additions are part of the XSI
21 extension; in other cases the standard developers decided that now was the time to migrate
22 these to the base standard.

23 The Single UNIX Specification programming environment provides a broad-based functional set
24 of interfaces to support the porting of existing UNIX applications and the development of new
25 applications. The environment also supports a rich set of tools for application development.

26 The majority of the obsolescent material from the existing POSIX.1 and POSIX.2 standards, and
27 material marked LEGACY from The Open Group’s Base specifications, has been removed in this
28 revision. New members of the Legacy Option Group have been added, reflecting the advance in
29 understanding of what is required.

30 The following IEEE standards have been added to the base documents in this revision:

- 31 • IEEE Std 1003.1d-1999
- 32 • IEEE Std 1003.1j-2000
- 33 • IEEE Std 1003.1q-2000
- 34 • IEEE P1003.1a draft standard
- 35 • IEEE Std 1003.2d-1994
- 36 • IEEE P1003.2b draft standard
- 37 • Selected parts of IEEE Std 1003.1g-2000

38 Only selected parts of IEEE Std 1003.1g-2000 were included. This was because there is much
39 duplication between the XNS, Issue 5.2 specification (another base document) and the material
40 from IEEE Std 1003.1g-2000, the former document being aligned with the latest networking
41 specifications for IPv6. Only the following sections of IEEE Std 1003.1g-2000 were considered for
42 inclusion:

- 43 • General terms related to sockets (Section 2.2.2)
- 44 • Socket concepts (Sections 5.1 through 5.3 inclusive)
- 45 • The *pselect()* function (Sections 6.2.2.1 and 6.2.3)
- 46 • The `<sys/select.h>` header (Section 6.2)

47 The following were requirements on IEEE Std 1003.1-2001:

- 48 • Backward-compatibility

49 It was agreed that there should be no breakage of functionality in the existing base
50 documents. This requirement was tempered by changes introduced due to interpretations
51 and corrigenda on the base documents, and any changes introduced in the
52 ISO/IEC 9899:1999 standard (C Language).

- 53 • Architecture and n-bit neutral

54 The common standard should not make any implicit assumptions about the system
55 architecture or size of data types; for example, previously some 32-bit implicit assumptions
56 had crept into the standards.

- 57 • Extensibility

58 It should be possible to extend the common standard without breaking backwards-
59 compatibility. For example, the name space should be reserved and structured to avoid
60 duplication of names between the standard and extensions to it.

61 **POSIX.1 and the ISO C Standard**

62 Previous revisions of POSIX.1 built upon the ISO C standard by reference only. This revision
63 takes a different approach.

64 The standard developers believed it essential for a programmer to have a single complete
65 reference place, but recognized that deference to the formal standard had to be addressed for the
66 duplicate interface definitions between the ISO C standard and the Single UNIX Specification.

67 It was agreed that where an interface has a version in the ISO C standard, the DESCRIPTION
68 section should describe the relationship to the ISO C standard and markings should be added as
69 appropriate to show where the ISO C standard has been extended in the text.

70 A block of text was added to the start of each affected reference page stating whether the page is
71 aligned with the ISO C standard or extended. Each page was parsed for additions beyond the
72 ISO C standard (that is, including both POSIX and UNIX extensions), and these extensions are
73 marked as CX extensions (for C Extensions).

74 **FIPS Requirements**

75 The Federal Information Processing Standards (FIPS) are a series of U.S. government
76 procurement standards managed and maintained on behalf of the U.S. Department of
77 Commerce by the National Institute of Standards and Technology (NIST).

78 The following restrictions have been made in this version of IEEE Std 1003.1 in order to align
79 with FIPS 151-2 requirements:

- 80 • The implementation supports `_POSIX_CHOWN_RESTRICTED`.
- 81 • The limit `{NGROUPS_MAX}` is now greater than or equal to 8.
- 82 • The implementation supports the setting of the group ID of a file (when it is created) to that
83 of the parent directory.

- 84 • The implementation supports `_POSIX_SAVED_IDS`.
- 85 • The implementation supports `_POSIX_VDISABLE`.
- 86 • The implementation supports `_POSIX_JOB_CONTROL`.
- 87 • The implementation supports `_POSIX_NO_TRUNC`.
- 88 • The `read()` function returns the number of bytes read when interrupted by a signal and does
89 not return `-1`.
- 90 • The `write()` function returns the number of bytes written when interrupted by a signal and
91 does not return `-1`.
- 92 • In the environment for the login shell, the environment variables `LOGNAME` and `HOME` are
93 defined and have the properties described in IEEE Std 1003.1-2001.
- 94 • The value of `{CHILD_MAX}` is now greater than or equal to 25.
- 95 • The value of `{OPEN_MAX}` is now greater than or equal to 20.
- 96 • The implementation supports the functionality associated with the symbols `CS7`, `CS8`,
97 `CSTOPB`, `PARODD`, and `PARENB` defined in `<termios.h>`.

98 A.1.2 Conformance

99 See Section A.2 (on page 9).

100 A.1.3 Normative References

101 There is no additional rationale provided for this section.

102 A.1.4 Terminology

103 The meanings specified in IEEE Std 1003.1-2001 for the words *shall*, *should*, and *may* are
104 mandated by ISO/IEC directives.

105 In the Rationale (Informative) volume of IEEE Std 1003.1-2001, the words *shall*, *should*, and *may*
106 are sometimes used to illustrate similar usages in IEEE Std 1003.1-2001. However, the rationale
107 itself does not specify anything regarding implementations or applications.

108 **conformance document**

109 As a practical matter, the conformance document is effectively part of the system
110 documentation. Conformance documents are distinguished by IEEE Std 1003.1-2001 so that
111 they can be referred to distinctly.

112 **implementation-defined**

113 This definition is analogous to that of the ISO C standard and, together with “undefined”
114 and “unspecified”, provides a range of specification of freedom allowed to the interface
115 implementor.

116 **may**

117 The use of *may* has been limited as much as possible, due both to confusion stemming from
118 its ordinary English meaning and to objections regarding the desirability of having as few
119 options as possible and those as clearly specified as possible.

120 The usage of *can* and *may* were selected to contrast optional application behavior (*can*)
121 against optional implementation behavior (*may*).

122 **shall**

123 Declarative sentences are sometimes used in IEEE Std 1003.1-2001 as if they included the
124 word *shall*, and facilities thus specified are no less required. For example, the two
125 statements:

126 1. The *foo()* function shall return zero.

127 2. The *foo()* function returns zero.

128 are meant to be exactly equivalent.

129 **should**

130 In IEEE Std 1003.1-2001, the word *should* does not usually apply to the implementation, but
131 rather to the application. Thus, the important words regarding implementations are *shall*,
132 which indicates requirements, and *may*, which indicates options.

133 **obsolescent**

134 The term “obsolescent” means “do not use this feature in new applications”. The
135 obsolescence concept is not an ideal solution, but was used as a method of increasing
136 consensus: many more objections would be heard from the user community if some of these
137 historical features were suddenly withdrawn without the grace period obsolescence
138 implies. The phrase “may be considered for withdrawal in future revisions” implies that the
139 result of that consideration might in fact keep those features indefinitely if the
140 predominance of applications do not migrate away from them quickly.

141 **legacy**

142 The term “legacy” was added for compatibility with the Single UNIX Specification. It
143 means “this feature is historic and optional; do not use this feature in new applications.
144 There are alternative interfaces that are more suitable.”. It is used exclusively for XSI
145 extensions, and includes facilities that were mandatory in previous versions of the base
146 document but are optional in this revision. This is a way to “sunset” the usage of certain
147 functions. Application writers should not rely on the existence of these facilities in new
148 applications, but should follow the migration path detailed in the APPLICATION USAGE
149 sections of the relevant pages.

150 The terms “legacy” and “obsolescent” are different: a feature marked LEGACY is not
151 recommended for new work and need not be present on an implementation (if the XSI
152 Legacy Option Group is not supported). A feature noted as obsolescent is supported by all
153 implementations, but may be removed in a future revision; new applications should not use
154 these features.

155 **system documentation**

156 The system documentation should normally describe the whole of the implementation,
157 including any extensions provided by the implementation. Such documents normally
158 contain information at least as detailed as the specifications in IEEE Std 1003.1-2001. Few
159 requirements are made on the system documentation, but the term is needed to avoid a
160 dangling pointer where the conformance document is permitted to point to the system
161 documentation.

162 **undefined**

163 See *implementation-defined*.

164 **unspecified**

165 See *implementation-defined*.

166 The definitions for “unspecified” and “undefined” appear nearly identical at first
167 examination, but are not. The term “unspecified” means that a conforming application may
168 deal with the unspecified behavior, and it should not care what the outcome is. The term

169 “undefined” says that a conforming application should not do it because no definition is
170 provided for what it does (and implicitly it would care what the outcome was if it tried it). It
171 is important to remember, however, that if the syntax permits the statement at all, it must
172 have some outcome in a real implementation.

173 Thus, the terms “undefined” and “unspecified” apply to the way the application should
174 think about the feature. In terms of the implementation, it is always “defined”—there is
175 always some result, even if it is an error. The implementation is free to choose the behavior
176 it prefers.

177 This also implies that an implementation, or another standard, could specify or define the
178 result in a useful fashion. The terms apply to IEEE Std 1003.1-2001 specifically.

179 The term “implementation-defined” implies requirements for documentation that are not
180 required for “undefined” (or “unspecified”). Where there is no need for a conforming
181 program to know the definition, the term “undefined” is used, even though
182 “implementation-defined” could also have been used in this context. There could be a
183 fourth term, specifying “this standard does not say what this does; it is acceptable to define
184 it in an implementation, but it does not need to be documented”, and undefined would then
185 be used very rarely for the few things for which any definition is not useful. In particular,
186 implementation-defined is used where it is believed that certain classes of application will
187 need to know such details to determine whether the application can be successfully ported
188 to the implementation. Such applications are not always strictly portable, but nevertheless
189 are common and useful; often the requirements met by the application cannot be met
190 without dealing with the issues implied by “implementation-defined”.

191 In many places IEEE Std 1003.1-2001 is silent about the behavior of some possible construct.
192 For example, a variable may be defined for a specified range of values and behaviors are
193 described for those values; nothing is said about what happens if the variable has any other
194 value. That kind of silence can imply an error in the standard, but it may also imply that the
195 standard was intentionally silent and that any behavior is permitted. There is a natural
196 tendency to infer that if the standard is silent, a behavior is prohibited. That is not the intent.
197 Silence is intended to be equivalent to the term “unspecified”.

198 The term “application” is not defined in IEEE Std 1003.1-2001; it is assumed to be a part of
199 general computer science terminology.

200 Three terms used within IEEE Std 1003.1-2001 overlap in meaning: “macro”, “symbolic name”,
201 and “symbolic constant”.

202 **macro**

203 This usually describes a C preprocessor symbol, the result of the **#define** operator, with or
204 without an argument. It may also be used to describe similar mechanisms in editors and
205 text processors.

206 **symbolic name**

207 This can also refer to a C preprocessor symbol (without arguments), but is also used to refer
208 to the names for characters in character sets. In addition, it is sometimes used to refer to
209 host names and even filenames.

210 **symbolic constant**

211 This also refers to a C preprocessor symbol (also without arguments).

212 In most cases, the difference in semantic content is negligible to nonexistent. Readers should not
213 attempt to read any meaning into the various usages of these terms.

214 **A.1.5 Portability**

215 To aid the identification of options within IEEE Std 1003.1-2001, a notation consisting of margin
 216 codes and shading is used. This is based on the notation used in previous revisions of The Open
 217 Group's Base specifications.

218 The benefit of this approach is a reduction in the number of *if* statements within the running
 219 text, that makes the text easier to read, and also an identification to the programmer that they
 220 need to ensure that their target platforms support the underlying options. For example, if
 221 functionality is marked with THR in the margin, it will be available on all systems supporting
 222 the Threads option, but may not be available on some others.

223 **A.1.5.1 Codes**

224 This section includes codes for options defined in the Base Definitions volume of
 225 IEEE Std 1003.1-2001, Section 2.1.6, Options, and the following additional codes for other
 226 purposes:

227 **CX** This margin code is used to denote extensions beyond the ISO C standard. For
 228 interfaces that are duplicated between IEEE Std 1003.1-2001 and the ISO C standard, a
 229 CX introduction block describes the nature of the duplication, with any extensions
 230 appropriately CX marked and shaded.

231 Where an interface is added to an ISO C standard header, within the header the
 232 interface has an appropriate margin marker and shading (for example, CX, XSI, TSF,
 233 and so on) and the same marking appears on the reference page in the SYNOPSIS
 234 section. This enables a programmer to easily identify that the interface is extending an
 235 ISO C standard header.

236 **MX** This margin code is used to denote IEC 60559: 1989 standard floating-point extensions.

237 **OB** This margin code is used to denote obsolescent behavior and thus flag a possible future
 238 applications portability warning.

239 **OH** The Single UNIX Specification has historically tried to reduce the number of headers an
 240 application has had to include when using a particular interface. Sometimes this was
 241 fewer than the base standard, and hence a notation is used to flag which headers are
 242 optional if you are using a system supporting the XSI extension.

243 **XSI** This code is used to denote interfaces and facilities within interfaces only required on
 244 systems supporting the XSI extension. This is introduced to support the Single UNIX
 245 Specification.

246 **XSR** This code is used to denote interfaces and facilities within interfaces only required on
 247 systems supporting STREAMS. This is introduced to support the Single UNIX
 248 Specification, although it is defined in a way so that it can stand alone from the XSI
 249 notation.

250 **A.1.5.2 Margin Code Notation**

251 Since some features may depend on one or more options, or require more than one option, a
 252 notation is used. Where a feature requires support of a single option, a single margin code will
 253 occur in the margin. If it depends on two options and both are required, then the codes will
 254 appear with a <space> separator. If either of two options are required, then a logical OR is
 255 denoted using the ' | ' symbol. If more than two codes are used, a special notation is used.

256 A.2 Conformance

257 The terms “profile” and “profiling” are used throughout this section.

258 A profile of a standard or standards is a codified set of option selections, such that by being
259 conformant to a profile, particular classes of users are specifically supported.

260 These conformance definitions are descended from those in the ISO POSIX-1: 1996 standard, but
261 with changes for the following:

- 262 • The addition of profiling options, allowing larger profiles of options such as the XSI
263 extension used by the Single UNIX Specification. In effect, it has profiled itself (that is,
264 created a self-profile).
- 265 • The addition of a definition of subprofiling considerations, to allow smaller profiles of
266 options.
- 267 • The addition of a hierarchy of super-options for XSI; these were formerly known as “Feature
268 Groups” in the System Interfaces and Headers, Issue 5 specification.
- 269 • Options from the ISO POSIX-2: 1993 standard are also now included, as IEEE Std 1003.1-2001
270 merges the functionality from it.

271 A.2.1 Implementation Conformance

272 These definitions allow application developers to know what to depend on in an
273 implementation.

274 There is no definition of a “strictly conforming implementation”; that would be an
275 implementation that provides *only* those facilities specified by POSIX.1 with no extensions
276 whatsoever. This is because no actual operating system implementation can exist without
277 system administration and initialization facilities that are beyond the scope of POSIX.1.

278 A.2.1.1 Requirements

279 The word “support” is used in certain instances, rather than “provide”, in order to allow an
280 implementation that has no resident software development facilities, but that supports the
281 execution of a *Strictly Conforming POSIX.1 Application*, to be a *conforming implementation*.

282 A.2.1.2 Documentation

283 The conformance documentation is required to use the same numbering scheme as POSIX.1 for
284 purposes of cross-referencing. All options that an implementation chooses are reflected in
285 <limits.h> and <unistd.h>.

286 Note that the use of “may” in terms of where conformance documents record where
287 implementations may vary, implies that it is not required to describe those features identified as
288 undefined or unspecified.

289 Other aspects of systems must be evaluated by purchasers for suitability. Many systems
290 incorporate buffering facilities, maintaining updated data in volatile storage and transferring
291 such updates to non-volatile storage asynchronously. Various exception conditions, such as a
292 power failure or a system crash, can cause this data to be lost. The data may be associated with a
293 file that is still open, with one that has been closed, with a directory, or with any other internal
294 system data structures associated with permanent storage. This data can be lost, in whole or
295 part, so that only careful inspection of file contents could determine that an update did not
296 occur.

297 Also, interrelated file activities, where multiple files and/or directories are updated, or where
298 space is allocated or released in the file system structures, can leave inconsistencies in the
299 relationship between data in the various files and directories, or in the file system itself. Such
300 inconsistencies can break applications that expect updates to occur in a specific sequence, so that
301 updates in one place correspond with related updates in another place.

302 For example, if a user creates a file, places information in the file, and then records this action in
303 another file, a system or power failure at this point followed by restart may result in a state in
304 which the record of the action is permanently recorded, but the file created (or some of its
305 information) has been lost. The consequences of this to the user may be undesirable. For a user
306 on such a system, the only safe action may be to require the system administrator to have a
307 policy that requires, after any system or power failure, that the entire file system must be
308 restored from the most recent backup copy (causing all intervening work to be lost).

309 The characteristics of each implementation will vary in this respect and may or may not meet the
310 requirements of a given application or user. Enforcement of such requirements is beyond the
311 scope of POSIX.1. It is up to the purchaser to determine what facilities are provided in an
312 implementation that affect the exposure to possible data or sequence loss, and also what
313 underlying implementation techniques and/or facilities are provided that reduce or limit such
314 loss or its consequences.

315 A.2.1.3 *POSIX Conformance*

316 This really means conformance to the base standard; however, since this revision includes the
317 core material of the Single UNIX Specification, the standard developers decided that it was
318 appropriate to segment the conformance requirements into two, the former for the base
319 standard, and the latter for the Single UNIX Specification.

320 Within POSIX.1 there are some symbolic constants that, if defined, indicate that a certain option
321 is enabled. Other symbolic constants exist in POSIX.1 for other reasons.

322 As part of the revision some alignment has occurred of the options with the FIPS 151-2 profile on
323 the POSIX.1-1990 standard. The following options from the POSIX.1-1990 standard are now
324 mandatory:

- 325 • `_POSIX_JOB_CONTROL`
- 326 • `_POSIX_SAVED_IDS`
- 327 • `_POSIX_VDISABLE`

328 A POSIX-conformant system may support the XSI extensions of the Single UNIX Specification.
329 This was intentional since the standard developers intend them to be upwards-compatible, so
330 that a system conforming to the Single UNIX Specification can also conform to the base standard
331 at the same time.

332 A.2.1.4 *XSI Conformance*

333 This section is added since the revision merges in the base volumes of the Single UNIX
334 Specification.

335 XSI conformance can be thought of as a profile, selecting certain options from
336 IEEE Std 1003.1-2001.

337 A.2.1.5 Option Groups

338 The concept of “Option Groups” is introduced to IEEE Std 1003.1-2001 to allow collections of
 339 related functions or options to be grouped together. This has been used as follows: the “XSI
 340 Option Groups” have been created to allow super-options, collections of underlying options and
 341 related functions, to be collectively supported by XSI-conforming systems. These reflect the
 342 “Feature Groups” from the System Interfaces and Headers, Issue 5 specification.

343 The standard developers considered the matter of subprofiling and decided it was better to
 344 include an enabling mechanism rather than detailed normative requirements. A set of
 345 subprofiling options was developed and included later in this volume of IEEE Std 1003.1-2001 as
 346 an informative illustration.

347 **Subprofiling Considerations**

348 The goal of not simultaneously fixing maximums and minimums was to allow implementations
 349 of the base standard or standards to support multiple profiles without conflict.

350 The following summarizes the rules for the limit types:

Limit Type	Fixed Value	Minimum Acceptable Value	Maximum Acceptable Value
Standard Profile	X_s $X_p == X_s$ (No change)	Y_s $Y_p \geq Y_s$ (May increase the limit)	Z_s $Z_p \leq Z_s$ (May decrease the limit)

356 The intent is that ranges specified by limits in profiles be entirely contained within the
 357 corresponding ranges of the base standard or standards being profiled, and that the unlimited
 358 end of a range in a base standard must remain unlimited in any profile of that standard.

359 Thus, the fixed `_POSIX_*` limits are constants and must not be changed by a profile. The
 360 variable counterparts (typically without the leading `_POSIX_`) can be changed but still remain
 361 semantically the same; that is, they still allow implementation values to vary as long as they
 362 meet the requirements for that value (be it a minimum or maximum).

363 Where a profile does not provide a feature upon which a limit is based, the limit is not relevant.
 364 Applications written to that profile should be written to operate independently of the value of
 365 the limit.

366 An example which has previously allowed implementations to support both the base standard
 367 and two other profiles in a compatible manner follows:

```
368     Base standard (POSIX.1-1996): _POSIX_CHILD_MAX 6
369     Base standard: CHILD_MAX    minimum maximum _POSIX_CHILD_MAX
370     FIPS profile/SUSv2 CHILD_MAX    25 (minimum maximum)
```

371 Another example:

```
372     Base standard (POSIX.1-1996): _POSIX_NGROUPS_MAX 0
373     Base standard: NGROUPS_MAX    minimum maximum _POSIX_NGROUP_MAX
374     FIPS profile/SUSv2 NGROUPS_MAX    8
```

375 A profile may lower a minimum maximum below the equivalent `_POSIX` value:

```
376     Base standard: _POSIX_foo_MAX    Z
377     Base standard: foo_MAX    _POSIX_foo_MAX
378     profile standard : foo_MAX    X (X can be less than, equal to,
379                                     or greater than _POSIX_foo_MAX)
```

380 In this case an implementation conforming to the profile may not conform to the base standard,
381 but an implementation to the base standard will conform to the profile.

382 A.2.1.6 Options

383 The final subsections within *Implementation Conformance* list the core options within
384 IEEE Std 1003.1-2001. This includes both options for the System Interfaces volume of
385 IEEE Std 1003.1-2001 and the Shell and Utilities volume of IEEE Std 1003.1-2001.

386 A.2.2 Application Conformance

387 These definitions guide users or adaptors of applications in determining on which
388 implementations an application will run and how much adaptation would be required to make
389 it run on others. These definitions are modeled after related ones in the ISO C standard.

390 POSIX.1 occasionally uses the expressions “portable application” or “conforming application”.
391 As they are used, these are synonyms for any of these terms. The differences between the classes
392 of application conformance relate to the requirements for other standards, the options supported
393 (such as the XSI extension) or, in the case of the Conforming POSIX.1 Application Using
394 Extensions, to implementation extensions. When one of the less explicit expressions is used, it
395 should be apparent from the context of the discussion which of the more explicit names is
396 appropriate

397 A.2.2.1 Strictly Conforming POSIX Application

398 This definition is analogous to that of an ISO C standard “conforming program”.

399 The major difference between a Strictly Conforming POSIX Application and an ISO C standard
400 strictly conforming program is that the latter is not allowed to use features of POSIX that are not
401 in the ISO C standard.

402 A.2.2.2 Conforming POSIX Application

403 Examples of <National Bodies> include ANSI, BSI, and AFNOR.

404 A.2.2.3 Conforming POSIX Application Using Extensions

405 Due to possible requirements for configuration or implementation characteristics in excess of the
406 specifications in <limits.h> or related to the hardware (such as array size or file space), not every
407 Conforming POSIX Application Using Extensions will run on every conforming
408 implementation.

409 A.2.2.4 Strictly Conforming XSI Application

410 This is intended to be upwards-compatible with the definition of a Strictly Conforming POSIX
411 Application, with the addition of the facilities and functionality included in the XSI extension.

412 A.2.2.5 Conforming XSI Application Using Extensions

413 Such applications may use extensions beyond the facilities defined by IEEE Std 1003.1-2001
414 including the XSI extension, but need to document the additional requirements.

415 **A.2.3 Language-Dependent Services for the C Programming Language**

416 POSIX.1 is, for historical reasons, both a specification of an operating system interface, shell and
417 utilities, and a C binding for that specification. Efforts had been previously undertaken to
418 generate a language-independent specification; however, that had failed, and the fact that the
419 ISO C standard is the *de facto* primary language on POSIX and the UNIX system makes this a
420 necessary and workable situation.

421 **A.2.4 Other Language-Related Specifications**

422 There is no additional rationale provided for this section.

423 **A.3 Definitions**

424 The definitions in this section are stated so that they can be used as exact substitutes for the
425 terms in text. They should not contain requirements or cross-references to sections within
426 IEEE Std 1003.1-2001; that is accomplished by using an informative note. In addition, the term
427 should not be included in its own definition. Where requirements or descriptions need to be
428 addressed but cannot be included in the definitions, due to not meeting the above criteria, these
429 occur in the General Concepts chapter.

430 In this revision, the definitions have been reworked extensively to meet style requirements and
431 to include terms from the base documents (see the Scope).

432 Many of these definitions are necessarily circular, and some of the terms (such as “process”) are
433 variants of basic computing science terms that are inherently hard to define. Where some
434 definitions are more conceptual and contain requirements, these appear in the General Concepts
435 chapter. Those listed in this section appear in an alphabetical glossary format of terms.

436 Some definitions must allow extension to cover terms or facilities that are not explicitly
437 mentioned in IEEE Std 1003.1-2001. For example, the definition of “Extended Security Controls”
438 permits implementations beyond those defined in IEEE Std 1003.1-2001.

439 Some terms in the following list of notes do not appear in IEEE Std 1003.1-2001; these are
440 marked prefixed with an asterisk (*). Many of them have been specifically excluded from
441 IEEE Std 1003.1-2001 because they concern system administration, implementation, or other
442 issues that are not specific to the programming interface. Those are marked with a reason, such
443 as “implementation-defined”.

444 **Appropriate Privileges**

445 One of the fundamental security problems with many historical UNIX systems has been that the
446 privilege mechanism is monolithic—a user has either no privileges or *all* privileges. Thus, a
447 successful “trojan horse” attack on a privileged process defeats all security provisions.
448 Therefore, POSIX.1 allows more granular privilege mechanisms to be defined. For many
449 historical implementations of the UNIX system, the presence of the term “appropriate
450 privileges” in POSIX.1 may be understood as a synonym for “superuser” (UID 0). However,
451 other systems have emerged where this is not the case and each discrete controllable action has
452 *appropriate privileges* associated with it. Because this mechanism is implementation-defined, it
453 must be described in the conformance document. Although that description affects several parts
454 of POSIX.1 where the term “appropriate privilege” is used, because the term “implementation-
455 defined” only appears here, the description of the entire mechanism and its effects on these
456 other sections belongs in this equivalent section of the conformance document. This is especially
457 convenient for implementations with a single mechanism that applies in all areas, since it only
458 needs to be described once.

459 **Byte**

460 The restriction that a byte is now exactly eight bits was a conscious decision by the standard
461 developers. It came about due to a combination of factors, primarily the use of the type `int8_t`
462 within the networking functions and the alignment with the ISO/IEC 9899:1999 standard, where
463 the `intN_t` types are now defined.

464 According to the ISO/IEC 9899:1999 standard:

- 465 • The `[u]intN_t` types must be two's complement with no padding bits and no illegal values.
- 466 • All types (apart from bit fields, which are not relevant here) must occupy an integral number
467 of bytes.
- 468 • If a type with width W occupies B bytes with C bits per byte (C is the value of `{CHAR_BIT}`),
469 then it has P padding bits where $P+W=B*C$.
- 470 • Therefore, for `int8_t` $P=0$, $W=8$. Since $B \geq 1$, $C \geq 8$, the only solution is $B=1$, $C=8$.

471 The standard developers also felt that this was not an undue restriction for the current state-of-
472 the-art for this version of IEEE Std 1003.1, but recognize that if industry trends continue, a wider
473 character type may be required in the future.

474 **Character**

475 The term “character” is used to mean a sequence of one or more bytes representing a single
476 graphic symbol. The deviation in the exact text of the ISO C standard definition for “byte” meets
477 the intent of the rationale of the ISO C standard also clears up the ambiguity raised by the term
478 “basic execution character set”. The octet-minimum requirement is a reflection of the
479 `{CHAR_BIT}` value.

480 **Clock Tick**

481 The ISO C standard defines a similar interval for use by the `clock()` function. There is no
482 requirement that these intervals be the same. In historical implementations these intervals are
483 different.

484 **Command**

485 The terms “command” and “utility” are related but have distinct meanings. Command is
486 defined as “a directive to a shell to perform a specific task”. The directive can be in the form of a
487 single utility name (for example, `ls`), or the directive can take the form of a compound command
488 (for example, `ls | grep name | pr`). A utility is a program that can be called by name
489 from a shell. Issuing only the name of the utility to a shell is the equivalent of a one-word
490 command. A utility may be invoked as a separate program that executes in a different process
491 than the command language interpreter, or it may be implemented as a part of the command
492 language interpreter. For example, the `echo` command (the directive to perform a specific task)
493 may be implemented such that the `echo` utility (the logic that performs the task of echoing) is in a
494 separate program; therefore, it is executed in a process that is different from the command
495 language interpreter. Conversely, the logic that performs the `echo` utility could be built into the
496 command language interpreter; therefore, it could execute in the same process as the command
497 language interpreter.

498 The terms “tool” and “application” can be thought of as being synonymous with “utility” from
499 the perspective of the operating system kernel. Tools, applications, and utilities historically have
500 run, typically, in processes above the kernel level. Tools and utilities historically have been a part
501 of the operating system non-kernel code and have performed system-related functions, such as
502 listing directory contents, checking file systems, repairing file systems, or extracting system

503 status information. Applications have not generally been a part of the operating system, and
504 they perform non-system-related functions, such as word processing, architectural design,
505 mechanical design, workstation publishing, or financial analysis. Utilities have most frequently
506 been provided by the operating system distributor, applications by third-party software
507 distributors, or by the users themselves. Nevertheless, IEEE Std 1003.1-2001 does not
508 differentiate between tools, utilities, and applications when it comes to receiving services from
509 the system, a shell, or the standard utilities. (For example, the *xargs* utility invokes another
510 utility; it would be of fairly limited usefulness if the users could not run their own applications
511 in place of the standard utilities.) Utilities are not applications in the sense that they are not
512 themselves subject to the restrictions of IEEE Std 1003.1-2001 or any other standard—there is no
513 requirement for *grep*, *stty*, or any of the utilities defined here to be any of the classes of
514 conforming applications.

515 **Column Positions**

516 In most 1-byte character sets, such as ASCII, the concept of column positions is identical to
517 character positions and to bytes. Therefore, it has been historically acceptable for some
518 implementations to describe line folding or tab stops or table column alignment in terms of bytes
519 or character positions. Other character sets pose complications, as they can have internal
520 representations longer than one octet and they can have display characters that have different
521 widths on the terminal screen or printer.

522 In IEEE Std 1003.1-2001 the term “column positions” has been defined to mean character—not
523 byte—positions in input files (such as “column position 7 of the FORTRAN input”). Output files
524 describe the column position in terms of the display width of the narrowest printable character
525 in the character set, adjusted to fit the characteristics of the output device. It is very possible that
526 *n* column positions will not be able to hold *n* characters in some character sets, unless all of those
527 characters are of the narrowest width. It is assumed that the implementation is aware of the
528 width of the various characters, deriving this information from the value of *LC_CTYPE*, and thus
529 can determine how many column positions to allot for each character in those utilities where it is
530 important.

531 The term “column position” was used instead of the more natural “column” because the latter is
532 frequently used in the different contexts of columns of figures, columns of table values, and so
533 on. Wherever confusion might result, these latter types of columns are referred to as “text
534 columns”.

535 **Controlling Terminal**

536 The question of which of possibly several special files referring to the terminal is meant is not
537 addressed in POSIX.1. The filename */dev/tty* is a synonym for the controlling terminal associated
538 with a process.

539 **Device Number***

540 The concept is handled in *stat()* as *ID of device*.

541 Direct I/O

542 Historically, direct I/O refers to the system bypassing intermediate buffering, but may be
543 extended to cover implementation-defined optimizations.

544 Directory

545 The format of the directory file is implementation-defined and differs radically between
546 System V and 4.3 BSD. However, routines (derived from 4.3 BSD) for accessing directories and
547 certain constraints on the format of the information returned by those routines are described in
548 the `<dirent.h>` header.

549 Directory Entry

550 Throughout IEEE Std 1003.1-2001, the term “link” is used (about the `link()` function, for
551 example) in describing the objects that point to files from directories.

552 Display

553 The Shell and Utilities volume of IEEE Std 1003.1-2001 assigns precise requirements for the
554 terms “display” and “write”. Some historical systems have chosen to implement certain utilities
555 without using the traditional file descriptor model. For example, the *vi* editor might employ
556 direct screen memory updates on a personal computer, rather than a `write()` system call. An
557 instance of user prompting might appear in a dialog box, rather than with standard error. When
558 the Shell and Utilities volume of IEEE Std 1003.1-2001 uses the term “display”, the method of
559 outputting to the terminal is unspecified; many historical implementations use *termcap* or
560 *terminfo*, but this is not a requirement. The term “write” is used when the Shell and Utilities
561 volume of IEEE Std 1003.1-2001 mandates that a file descriptor be used and that the output can
562 be redirected. However, it is assumed that when the writing is directly to the terminal (it has not
563 been redirected elsewhere), there is no practical way for a user or test suite to determine whether
564 a file descriptor is being used. Therefore, the use of a file descriptor is mandated only for the
565 redirection case and the implementation is free to use any method when the output is not
566 redirected. The verb *write* is used almost exclusively, with the very few exceptions of those
567 utilities where output redirection need not be supported: *tabs*, *talk*, *tput*, and *vi*.

568 Dot

569 The symbolic name *dot* is carefully used in POSIX.1 to distinguish the working directory
570 filename from a period or a decimal point.

571 Dot-Dot

572 Historical implementations permit the use of these filenames without their special meanings.
573 Such use precludes any meaningful use of these filenames by a Conforming POSIX.1
574 Application. Therefore, such use is considered an extension, the use of which makes an
575 implementation non-conforming; see also Section A.4.11 (on page 38).

576 Epoch

577 Historically, the origin of UNIX system time was referred to as “00:00:00 GMT, January 1, 1970”.
578 Greenwich Mean Time is actually not a term acknowledged by the international standards
579 community; therefore, this term, “Epoch”, is used to abbreviate the reference to the actual
580 standard, Coordinated Universal Time.

581 FIFO Special File

582 See **Pipe** (on page 24).

583 File

584 It is permissible for an implementation-defined file type to be non-readable or non-writable.

585 File Classes

586 These classes correspond to the historical sets of permission bits. The classes are general to
587 allow implementations flexibility in expanding the access mechanism for more stringent security
588 environments. Note that a process is in one and only one class, so there is no ambiguity.

589 Filename

590 At the present time, the primary responsibility for truncating filenames containing multi-byte
591 characters must reside with the application. Some industry groups involved in
592 internationalization believe that in the future the responsibility must reside with the kernel. For
593 the moment, a clearer understanding of the implications of making the kernel responsible for
594 truncation of multi-byte filenames is needed.

595 Character-level truncation was not adopted because there is no support in POSIX.1 that advises
596 how the kernel distinguishes between single and multi-byte characters. Until that time, it must
597 be incumbent upon application writers to determine where multi-byte characters must be
598 truncated.

599 File System

600 Historically, the meaning of this term has been overloaded with two meanings: that of the
601 complete file hierarchy, and that of a mountable subset of that hierarchy; that is, a mounted file
602 system. POSIX.1 uses the term “file system” in the second sense, except that it is limited to the
603 scope of a process (and a process’ root directory). This usage also clarifies the domain in which a
604 file serial number is unique.

605 Graphic Character

606 This definition is made available for those definitions (in particular, *TZ*) which must exclude
607 control characters.

608 Group Database

609 See **User Database** (on page 32).

610 Group File*

611 Implementation-defined; see **User Database** (on page 32).

612 Historical Implementations*

613 This refers to previously existing implementations of programming interfaces and operating
614 systems that are related to the interface specified by POSIX.1.

615 Hosted Implementation*

616 This refers to a POSIX.1 implementation that is accomplished through interfaces from the
617 POSIX.1 services to some alternate form of operating system kernel services. Note that the line
618 between a hosted implementation and a native implementation is blurred, since most
619 implementations will provide some services directly from the kernel and others through some
620 indirect path. (For example, *fopen()* might use *open()*; or *mkfifo()* might use *mknod()*.) There is
621 no necessary relationship between the type of implementation and its correctness, performance,
622 and/or reliability.

623 Implementation*

624 This term is generally used instead of its synonym, “system”, to emphasize the consequences of
625 decisions to be made by system implementors. Perhaps if no options or extensions to POSIX.1
626 were allowed, this usage would not have occurred.

627 The term “specific implementation” is sometimes used as a synonym for “implementation”.
628 This should not be interpreted too narrowly; both terms can represent a relatively broad group
629 of systems. For example, a hardware vendor could market a very wide selection of systems that
630 all used the same instruction set, with some systems desktop models and others large multi-user
631 minicomputers. This wide range would probably share a common POSIX.1 operating system,
632 allowing an application compiled for one to be used on any of the others; this is a [*specific*]
633 *implementation*. However, such a wide range of machines probably has some differences
634 between the models. Some may have different clock rates, different file systems, different
635 resource limits, different network connections, and so on, depending on their sizes or intended
636 usages. Even on two identical machines, the system administrators may configure them
637 differently. Each of these different systems is known by the term “a specific instance of a specific
638 implementation”. This term is only used in the portions of POSIX.1 dealing with runtime
639 queries: *sysconf()* and *pathconf()*.

640 Incomplete Pathname*

641 Absolute pathname has been adequately defined.

642 Job Control

643 In order to understand the job control facilities in POSIX.1 it is useful to understand how they
644 are used by a job control-cognizant shell to create the user interface effect of job control.

645 While the job control facilities supplied by POSIX.1 can, in theory, support different types of
646 interactive job control interfaces supplied by different types of shells, there was historically one
647 particular interface that was most common when the standard was originally developed
648 (provided by BSD C Shell). This discussion describes that interface as a means of illustrating
649 how the POSIX.1 job control facilities can be used.

650 Job control allows users to selectively stop (suspend) the execution of processes and continue
651 (resume) their execution at a later point. The user typically employs this facility via the
652 interactive interface jointly supplied by the terminal I/O driver and a command interpreter

653 (shell).

654 The user can launch jobs (command pipelines) in either the foreground or background. When
655 launched in the foreground, the shell waits for the job to complete before prompting for
656 additional commands. When launched in the background, the shell does not wait, but
657 immediately prompts for new commands.

658 If the user launches a job in the foreground and subsequently regrets this, the user can type the
659 suspend character (typically set to <control>-Z), which causes the foreground job to stop and the
660 shell to begin prompting for new commands. The stopped job can be continued by the user (via
661 special shell commands) either as a foreground job or as a background job. Background jobs can
662 also be moved into the foreground via shell commands.

663 If a background job attempts to access the login terminal (controlling terminal), it is stopped by
664 the terminal driver and the shell is notified, which, in turn, notifies the user. (Terminal access
665 includes *read()* and certain terminal control functions, and conditionally includes *write()*.) The
666 user can continue the stopped job in the foreground, thus allowing the terminal access to
667 succeed in an orderly fashion. After the terminal access succeeds, the user can optionally move
668 the job into the background via the suspend character and shell commands.

669 *Implementing Job Control Shells*

670 The interactive interface described previously can be accomplished using the POSIX.1 job
671 control facilities in the following way.

672 The key feature necessary to provide job control is a way to group processes into jobs. This
673 grouping is necessary in order to direct signals to a single job and also to identify which job is in
674 the foreground. (There is at most one job that is in the foreground on any controlling terminal at
675 a time.)

676 The concept of process groups is used to provide this grouping. The shell places each job in a
677 separate process group via the *setpgid()* function. To do this, the *setpgid()* function is invoked by
678 the shell for each process in the job. It is actually useful to invoke *setpgid()* twice for each
679 process: once in the child process, after calling *fork()* to create the process, but before calling one
680 of the *exec* family of functions to begin execution of the program, and once in the parent shell
681 process, after calling *fork()* to create the child. The redundant invocation avoids a race condition
682 by ensuring that the child process is placed into the new process group before either the parent
683 or the child relies on this being the case. The process group ID for the job is selected by the shell
684 to be equal to the process ID of one of the processes in the job. Some shells choose to make one
685 process in the job be the parent of the other processes in the job (if any). Other shells (for
686 example, the C Shell) choose to make themselves the parent of all processes in the pipeline (job).
687 In order to support this latter case, the *setpgid()* function accepts a process group ID parameter
688 since the correct process group ID cannot be inherited from the shell. The shell itself is
689 considered to be a job and is the sole process in its own process group.

690 The shell also controls which job is currently in the foreground. A foreground and background
691 job differ in two ways: the shell waits for a foreground command to complete (or stop) before
692 continuing to read new commands, and the terminal I/O driver inhibits terminal access by
693 background jobs (causing the processes to stop). Thus, the shell must work cooperatively with
694 the terminal I/O driver and have a common understanding of which job is currently in the
695 foreground. It is the user who decides which command should be currently in the foreground,
696 and the user informs the shell via shell commands. The shell, in turn, informs the terminal I/O
697 driver via the *tcsetgrp()* function. This indicates to the terminal I/O driver the process group ID
698 of the foreground process group (job). When the current foreground job either stops or
699 terminates, the shell places itself in the foreground via *tcsetgrp()* before prompting for
700 additional commands. Note that when a job is created the new process group begins as a

701 background process group. It requires an explicit act of the shell via *tcsetpgrp()* to move a
702 process group (job) into the foreground.

703 When a process in a job stops or terminates, its parent (for example, the shell) receives
704 synchronous notification by calling the *waitpid()* function with the WUNTRACED flag set.
705 Asynchronous notification is also provided when the parent establishes a signal handler for
706 SIGCHLD and does not specify the SA_NOCLDSTOP flag. Usually all processes in a job stop as
707 a unit since the terminal I/O driver always sends job control stop signals to all processes in the
708 process group.

709 To continue a stopped job, the shell sends the SIGCONT signal to the process group of the job. In
710 addition, if the job is being continued in the foreground, the shell invokes *tcsetpgrp()* to place the
711 job in the foreground before sending SIGCONT. Otherwise, the shell leaves itself in the
712 foreground and reads additional commands.

713 There is additional flexibility in the POSIX.1 job control facilities that allows deviations from the
714 typical interface. Clearing the TOSTOP terminal flag allows background jobs to perform *write()*
715 functions without stopping. The same effect can be achieved on a per-process basis by having a
716 process set the signal action for SIGTTOU to SIG_IGN.

717 Note that the terms “job” and “process group” can be used interchangeably. A login session that
718 is not using the job control facilities can be thought of as a large collection of processes that are
719 all in the same job (process group). Such a login session may have a partial distinction between
720 foreground and background processes; that is, the shell may choose to wait for some processes
721 before continuing to read new commands and may not wait for other processes. However, the
722 terminal I/O driver will consider all these processes to be in the foreground since they are all
723 members of the same process group.

724 In addition to the basic job control operations already mentioned, a job control-cognizant shell
725 needs to perform the following actions.

726 When a foreground (not background) job stops, the shell must sample and remember the current
727 terminal settings so that it can restore them later when it continues the stopped job in the
728 foreground (via the *tcgetattr()* and *tcsetattr()* functions).

729 Because a shell itself can be spawned from a shell, it must take special action to ensure that
730 subshells interact well with their parent shells.

731 A subshell can be spawned to perform an interactive function (prompting the terminal for
732 commands) or a non-interactive function (reading commands from a file). When operating non-
733 interactively, the job control shell will refrain from performing the job control-specific actions
734 described above. It will behave as a shell that does not support job control. For example, all jobs
735 will be left in the same process group as the shell, which itself remains in the process group
736 established for it by its parent. This allows the shell and its children to be treated as a single job
737 by a parent shell, and they can be affected as a unit by terminal keyboard signals.

738 An interactive subshell can be spawned from another job control-cognizant shell in either the
739 foreground or background. (For example, from the C Shell, the user can execute the command,
740 *csH &*.) Before the subshell activates job control by calling *setpgid()* to place itself in its own
741 process group and *tcsetpgrp()* to place its new process group in the foreground, it needs to
742 ensure that it has already been placed in the foreground by its parent. (Otherwise, there could
743 be multiple job control shells that simultaneously attempt to control mediation of the terminal.)
744 To determine this, the shell retrieves its own process group via *getpgrp()* and the process group
745 of the current foreground job via *tcgetpgrp()*. If these are not equal, the shell sends SIGTTIN to
746 its own process group, causing itself to stop. When continued later by its parent, the shell
747 repeats the process group check. When the process groups finally match, the shell is in the
748 foreground and it can proceed to take control. After this point, the shell ignores all the job

749 control stop signals so that it does not inadvertently stop itself.

750 *Implementing Job Control Applications*

751 Most applications do not need to be aware of job control signals and operations; the intuitively
752 correct behavior happens by default. However, sometimes an application can inadvertently
753 interfere with normal job control processing, or an application may choose to overtly effect job
754 control in cooperation with normal shell procedures.

755 An application can inadvertently subvert job control processing by “blindly” altering the
756 handling of signals. A common application error is to learn how many signals the system
757 supports and to ignore or catch them all. Such an application makes the assumption that it does
758 not know what this signal is, but knows the right handling action for it. The system may
759 initialize the handling of job control stop signals so that they are being ignored. This allows
760 shells that do not support job control to inherit and propagate these settings and hence to be
761 immune to stop signals. A job control shell will set the handling to the default action and
762 propagate this, allowing processes to stop. In doing so, the job control shell is taking
763 responsibility for restarting the stopped applications. If an application wishes to catch the stop
764 signals itself, it should first determine their inherited handling states. If a stop signal is being
765 ignored, the application should continue to ignore it. This is directly analogous to the
766 recommended handling of SIGINT described in the referenced UNIX Programmer’s Manual.

767 If an application is reading the terminal and has disabled the interpretation of special characters
768 (by clearing the ISIG flag), the terminal I/O driver will not send SIGTSTP when the suspend
769 character is typed. Such an application can simulate the effect of the suspend character by
770 recognizing it and sending SIGTSTP to its process group as the terminal driver would have
771 done. Note that the signal is sent to the process group, not just to the application itself; this
772 ensures that other processes in the job also stop. (Note also that other processes in the job could
773 be children, siblings, or even ancestors.) Applications should not assume that the suspend
774 character is <control>-Z (or any particular value); they should retrieve the current setting at
775 startup.

776 *Implementing Job Control Systems*

777 The intent in adding 4.2 BSD-style job control functionality was to adopt the necessary 4.2 BSD
778 programmatic interface with only minimal changes to resolve syntactic or semantic conflicts
779 with System V or to close recognized security holes. The goal was to maximize the ease of
780 providing both conforming implementations and Conforming POSIX.1 Applications.

781 It is only useful for a process to be affected by job control signals if it is the descendant of a job
782 control shell. Otherwise, there will be nothing that continues the stopped process.

783 POSIX.1 does not specify how controlling terminal access is affected by a user logging out (that
784 is, by a controlling process terminating). 4.2 BSD uses the *vhangup()* function to prevent any
785 access to the controlling terminal through file descriptors opened prior to logout. System V does
786 not prevent controlling terminal access through file descriptors opened prior to logout (except
787 for the case of the special file, */dev/tty*). Some implementations choose to make processes
788 immune from job control after logout (that is, such processes are always treated as if in the
789 foreground); other implementations continue to enforce foreground/background checks after
790 logout. Therefore, a Conforming POSIX.1 Application should not attempt to access the
791 controlling terminal after logout since such access is unreliable. If an implementation chooses to
792 deny access to a controlling terminal after its controlling process exits, POSIX.1 requires a certain
793 type of behavior (see **Controlling Terminal** (on page 15)).

794 **Kernel***795 See **System Call*** (on page 31).796 **Library Routine***797 See **System Call*** (on page 31).798 **Logical Device***

799 Implementation-defined.

800 **Map**801 The definition of map is included to clarify the usage of mapped pages in the description of the
802 behavior of process memory locking.803 **Memory-Resident**804 The term “memory-resident” is historically understood to mean that the so-called resident
805 pages are actually present in the physical memory of the computer system and are immune from
806 swapping, paging, copy-on-write faults, and so on. This is the actual intent of
807 IEEE Std 1003.1-2001 in the process memory locking section for implementations where this is
808 logical. But for some implementations—primarily mainframes—actually locking pages into
809 primary storage is not advantageous to other system objectives, such as maximizing throughput.
810 For such implementations, memory locking is a “hint” to the implementation that the
811 application wishes to avoid situations that would cause long latencies in accessing memory.
812 Furthermore, there are other implementation-defined issues with minimizing memory access
813 latencies that “memory residency” does not address—such as MMU reload faults. The definition
814 attempts to accommodate various implementations while allowing conforming applications to
815 specify to the implementation that they want or need the best memory access times that the
816 implementation can provide.817 **Memory Object***818 The term “memory object” usually implies shared memory. If the object is the same as a
819 filename in the file system name space of the implementation, it is expected that the data written
820 into the memory object be preserved on disk. A memory object may also apply to a physical
821 device on an implementation. In this case, writes to the memory object are sent to the controller
822 for the device and reads result in control registers being returned.823 **Mount Point***824 The directory on which a “mounted file system” is mounted. This term, like *mount()* and
825 *umount()*, was not included because it was implementation-defined.826 **Mounted File System***827 See **File System** (on page 17).

828 **Name**

829 There are no explicit limits in IEEE Std 1003.1-2001 on the sizes of names, words (see the
830 definition of word in the Base Definitions volume of IEEE Std 1003.1-2001), lines, or other
831 objects. However, other implicit limits do apply: shell script lines produced by many of the
832 standard utilities cannot exceed {LINE_MAX} and the sum of exported variables comes under
833 the {ARG_MAX} limit. Historical shells dynamically allocate memory for names and words and
834 parse incoming lines a character at a time. Lines cannot have an arbitrary {LINE_MAX} limit
835 because of historical practice, such as makefiles, where *make* removes the <newline>s associated
836 with the commands for a target and presents the shell with one very long line. The text on
837 INPUT FILES in the Shell and Utilities volume of IEEE Std 1003.1-2001, Section 1.11, Utility
838 Description Defaults does allow a shell to run out of memory, but it cannot have arbitrary
839 programming limits.

840 **Native Implementation***

841 This refers to an implementation of POSIX.1 that interfaces directly to an operating system
842 kernel; see also *hosted implementation* and *cooperating implementation*. A similar concept is a
843 native UNIX system, which would be a kernel derived from one of the original UNIX system
844 products.

845 **Nice Value**

846 This definition is not intended to suggest that all processes in a system have priorities that are
847 comparable. Scheduling policy extensions, such as adding realtime priorities, make the notion of
848 a single underlying priority for all scheduling policies problematic. Some implementations may
849 implement the features related to *nice* to affect all processes on the system, others to affect just
850 the general time-sharing activities implied by IEEE Std 1003.1-2001, and others may have no
851 effect at all. Because of the use of “implementation-defined” in *nice* and *renice*, a wide range of
852 implementation strategies is possible.

853 **Open File Description**

854 An “open file description”, as it is currently named, describes how a file is being accessed. What
855 is currently called a “file descriptor” is actually just an identifier or “handle”; it does not actually
856 describe anything.

857 The following alternate names were discussed:

- 858 • For “open file description”:
859 “open instance”, “file access description”, “open file information”, and “file access
860 information”.
- 861 • For “file descriptor”:
862 “file handle”, “file number” (cf., *fileno()*). Some historical implementations use the term “file
863 table entry”.

864 **Orphaned Process Group**

865 Historical implementations have a concept of an orphaned process, which is a process whose
866 parent process has exited. When job control is in use, it is necessary to prevent processes from
867 being stopped in response to interactions with the terminal after they no longer are controlled by
868 a job control-cognizant program. Because signals generated by the terminal are sent to a process
869 group and not to individual processes, and because a signal may be provoked by a process that
870 is not orphaned, but sent to another process that is orphaned, it is necessary to define an
871 orphaned process group. The definition assumes that a process group will be manipulated as a

872 group and that the job control-cognizant process controlling the group is outside of the group
873 and is the parent of at least one process in the group (so that state changes may be reported via
874 *waitpid()*). Therefore, a group is considered to be controlled as long as at least one process in the
875 group has a parent that is outside of the process group, but within the session.

876 This definition of orphaned process groups ensures that a session leader's process group is
877 always considered to be orphaned, and thus it is prevented from stopping in response to
878 terminal signals.

879 **Page**

880 The term “page” is defined to support the description of the behavior of memory mapping for
881 shared memory and memory mapped files, and the description of the behavior of process
882 memory locking. It is not intended to imply that shared memory/file mapping and memory
883 locking are applicable only to “paged” architectures. For the purposes of IEEE Std 1003.1-2001,
884 whatever the granularity on which an architecture supports mapping or locking, this is
885 considered to be a “page”. If an architecture cannot support the memory mapping or locking
886 functions specified by IEEE Std 1003.1-2001 on any granularity, then these options will not be
887 implemented on the architecture.

888 **Passwd File***

889 Implementation-defined; see **User Database** (on page 32).

890 **Parent Directory**

891 There may be more than one directory entry pointing to a given directory in some
892 implementations. The wording here identifies that exactly one of those is the parent directory. In
893 pathname resolution, dot-dot is identified as the way that the unique directory is identified.
894 (That is, the parent directory is the one to which dot-dot points.) In the case of a remote file
895 system, if the same file system is mounted several times, it would appear as if they were distinct
896 file systems (with interesting synchronization properties).

897 **Pipe**

898 It proved convenient to define a pipe as a special case of a FIFO, even though historically the
899 latter was not introduced until System III and does not exist at all in 4.3 BSD.

900 **Portable Filename Character Set**

901 The encoding of this character set is not specified—specifically, ASCII is not required. But the
902 implementation must provide a unique character code for each of the printable graphics
903 specified by POSIX.1; see also Section A.4.6 (on page 34).

904 Situations where characters beyond the portable filename character set (or historically ASCII or
905 the ISO/IEC 646:1991 standard) would be used (in a context where the portable filename
906 character set or the ISO/IEC 646:1991 standard is required by POSIX.1) are expected to be
907 common. Although such a situation renders the use technically non-compliant, mutual
908 agreement among the users of an extended character set will make such use portable between
909 those users. Such a mutual agreement could be formalized as an optional extension to POSIX.1.
910 (Making it required would eliminate too many possible systems, as even those systems using the
911 ISO/IEC 646:1991 standard as a base character set extend their character sets for Western
912 Europe and the rest of the world in different ways.)

913 Nothing in POSIX.1 is intended to preclude the use of extended characters where interchange is
914 not required or where mutual agreement is obtained. It has been suggested that in several places

915 “should” be used instead of “shall”. Because (in the worst case) use of any character beyond the
916 portable filename character set would render the program or data not portable to all possible
917 systems, no extensions are permitted in this context.

918 **Regular File**

919 POSIX.1 does not intend to preclude the addition of structuring data (for example, record
920 lengths) in the file, as long as such data is not visible to an application that uses the features
921 described in POSIX.1.

922 **Root Directory**

923 This definition permits the operation of *chroot()*, even though that function is not in POSIX.1; see
924 also Section A.4.5 (on page 34).

925 **Root File System***

926 Implementation-defined.

927 **Root of a File System***

928 Implementation-defined; see **Mount Point*** (on page 22).

929 **Signal**

930 The definition implies a double meaning for the term. Although a signal is an event, common
931 usage implies that a signal is an identifier of the class of event.

932 **Superuser***

933 This concept, with great historical significance to UNIX system users, has been replaced with the
934 notion of appropriate privileges.

935 **Supplementary Group ID**

936 The POSIX.1-1990 standard is inconsistent in its treatment of supplementary groups. The
937 definition of supplementary group ID explicitly permits the effective group ID to be included in
938 the set, but wording in the description of the *setuid()* and *setgid()* functions states: “Any
939 supplementary group IDs of the calling process remain unchanged by these function calls”. In
940 the case of *setgid()* this contradicts that definition. In addition, some felt that the unspecified
941 behavior in the definition of supplementary group IDs adds unnecessary portability problems.
942 The standard developers considered several solutions to this problem:

- 943 1. Reword the description of *setgid()* to permit it to change the supplementary group IDs to
944 reflect the new effective group ID. A problem with this is that it adds more “may”s to the
945 wording and does not address the portability problems of this optional behavior.
- 946 2. Mandate the inclusion of the effective group ID in the supplementary set (giving
947 {NGROUPS_MAX} a minimum value of 1). This is the behavior of 4.4 BSD. In that system,
948 the effective group ID is the first element of the array of supplementary group IDs (there is
949 no separate copy stored, and changes to the effective group ID are made only in the
950 supplementary group set). By convention, the initial value of the effective group ID is
951 duplicated elsewhere in the array so that the initial value is not lost when executing a set-
952 group-ID program.
- 953 3. Change the definition of supplementary group ID to exclude the effective group ID and
954 specify that the effective group ID does not change the set of supplementary group IDs.

- 955 This is the behavior of 4.2 BSD, 4.3 BSD, and System V Release 4.
- 956 4. Change the definition of supplementary group ID to exclude the effective group ID, and
 957 require that *getgroups()* return the union of the effective group ID and the supplementary
 958 group IDs.
- 959 5. Change the definition of {NGROUPS_MAX} to be one more than the number of
 960 supplementary group IDs, so it continues to be the number of values returned by
 961 *getgroups()* and existing applications continue to work. This alternative is effectively the
 962 same as the second (and might actually have the same implementation).

963 The standard developers decided to permit either 2 or 3. The effective group ID is orthogonal to
 964 the set of supplementary group IDs, and it is implementation-defined whether *getgroups()*
 965 returns this. If the effective group ID is returned with the set of supplementary group IDs, then
 966 all changes to the effective group ID affect the supplementary group set returned by *getgroups()*.
 967 It is permissible to eliminate duplicates from the list returned by *getgroups()*. However, if a
 968 group ID is contained in the set of supplementary group IDs, setting the group ID to that value
 969 and then to a different value should not remove that value from the supplementary group IDs.

970 The definition of supplementary group IDs has been changed to not include the effective group
 971 ID. This simplifies permanent rationale and makes the relevant functions easier to understand.
 972 The *getgroups()* function has been modified so that it can, on an implementation-defined basis,
 973 return the effective group ID. By making this change, functions that modify the effective group
 974 ID do not need to discuss adding to the supplementary group list; the only view into the
 975 supplementary group list that the application writer has is through the *getgroups()* function.

976 Symbolic Link

977 Many implementations associate no attributes, including ownership with symbolic links.
 978 Security experts encouraged consideration for defining these attributes as optional.
 979 Consideration was given to changing *utime()* to allow modification of the times for a symbolic
 980 link, or as an alternative adding an *lutime()* interface. Modifications to *chown()* were also
 981 considered: allow changing symbolic link ownership or alternatively adding *lchown()*. As a
 982 result of alignment with the Single UNIX Specification, the *lchown()* function is included as part
 983 of the XSI extension and XSI-conformant systems may support an owner and a group associated
 984 with a symbolic link. There was no consensus to define further attributes for symbolic links, and
 985 for systems not supporting the XSI extension only the file type bits in the *st_mode* member and
 986 the *st_size* member of the *stat* structure are required to be applicable to symbolic links.

987 Historical implementations were followed when determining which interfaces should apply to
 988 symbolic links. Interfaces that historically followed symbolic links include *chmod()*, *link()*, and
 989 *utime()*. Interfaces that historically do not follow symbolic links include *chown()*, *lstat()*,
 990 *readlink()*, *rename()*, *remove()*, *rmdir()*, and *unlink()*. IEEE Std 1003.1-2001 deviates from
 991 historical practice only in the case of *chown()*. Because there is no requirement for systems not
 992 supporting the XSI extension that there is an association of ownership with symbolic links, there
 993 was no interface in the base standard to change ownership. In addition, other implementations
 994 of symbolic links have modified *chown()* to follow symbolic links.

995 In the case of symbolic links, IEEE Std 1003.1-2001 states that a trailing slash is considered to be
 996 the final component of a pathname rather than the pathname component that preceded it. This is
 997 the behavior of historical implementations. For example, for */a/b* and */a/b/*, if */a/b* is a symbolic
 998 link to a directory, then */a/b* refers to the symbolic link, and */a/b/* is the same as */a/b/.*, which is the
 999 directory to which the symbolic link points.

1000 For multi-level security purposes, it is possible to have the link read mode govern permission for
 1001 the *readlink()* function. It is also possible that the read permissions of the directory containing

1002 the link be used for this purpose. Implementations may choose to use either of these methods;
1003 however, this is not current practice and neither method is specified.

1004 Several reasons were advanced for requiring that when a symbolic link is used as the source
1005 argument to the *link()* function, the resulting link will apply to the file named by the contents of
1006 the symbolic link rather than to the symbolic link itself. This is the case in historical
1007 implementations. This action was preferred, as it supported the traditional idea of persistence
1008 with respect to the target of a hard link. This decision is appropriate in light of a previous
1009 decision not to require association of attributes with symbolic links, thereby allowing
1010 implementations which do not use inodes. Opposition centered on the lack of symmetry on the
1011 part of the *link()* and *unlink()* function pair with respect to symbolic links.

1012 Because a symbolic link and its referenced object coexist in the file system name space, confusion
1013 can arise in distinguishing between the link itself and the referenced object. Historically, utilities
1014 and system calls have adopted their own link following conventions in a somewhat *ad hoc*
1015 fashion. Rules for a uniform approach are outlined here, although historical practice has been
1016 adhered to as much as was possible. To promote consistent system use, user-written utilities are
1017 encouraged to follow these same rules.

1018 Symbolic links are handled either by operating on the link itself, or by operating on the object
1019 referenced by the link. In the latter case, an application or system call is said to “follow” the link.
1020 Symbolic links may reference other symbolic links, in which case links are dereferenced until an
1021 object that is not a symbolic link is found, a symbolic link that references a file that does not exist
1022 is found, or a loop is detected. (Current implementations do not detect loops, but have a limit on
1023 the number of symbolic links that they will dereference before declaring it an error.)

1024 There are four domains for which default symbolic link policy is established in a system. In
1025 almost all cases, there are utility options that override this default behavior. The four domains
1026 are as follows:

- 1027 1. Symbolic links specified to system calls that take filename arguments
- 1028 2. Symbolic links specified as command line filename arguments to utilities that are not
1029 performing a traversal of a file hierarchy
- 1030 3. Symbolic links referencing files not of type directory, specified to utilities that are
1031 performing a traversal of a file hierarchy
- 1032 4. Symbolic links referencing files of type directory, specified to utilities that are performing a
1033 traversal of a file hierarchy

1034 *First Domain*

1035 The first domain is considered in earlier rationale.

1036 *Second Domain*

1037 The reason this category is restricted to utilities that are not traversing the file hierarchy is that
1038 some standard utilities take an option that specifies a hierarchical traversal, but by default
1039 operate on the arguments themselves. Generally, users specifying the option for a file hierarchy
1040 traversal wish to operate on a single, physical hierarchy, and therefore symbolic links, which
1041 may reference files outside of the hierarchy, are ignored. For example, *chown owner file* is a
1042 different operation from the same command with the **-R** option specified. In this example, the
1043 behavior of the command *chown owner file* is described here, while the behavior of the command
1044 *chown -R owner file* is described in the third and fourth domains.

1045 The general rule is that the utilities in this category follow symbolic links named as arguments.

1046 Exceptions in the second domain are:

- 1047 • The *mv* and *rm* utilities do not follow symbolic links named as arguments, but respectively
1048 attempt to rename or delete them.
- 1049 • The *ls* utility is also an exception to this rule. For compatibility with historical systems, when
1050 the **-R** option is not specified, the *ls* utility follows symbolic links named as arguments if the
1051 **-L** option is specified or if the **-F**, **-d**, or **-l** options are not specified. (If the **-L** option is
1052 specified, *ls* always follows symbolic links; it is the only utility where the **-L** option affects its
1053 behavior even though a tree walk is not being performed.)

1054 All other standard utilities, when not traversing a file hierarchy, always follow symbolic links
1055 named as arguments.

1056 Historical practice is that the **-h** option is specified if standard utilities are to act upon symbolic
1057 links instead of upon their targets. Examples of commands that have historically had a **-h** option
1058 for this purpose are the *chgrp*, *chown*, *file*, and *test* utilities.

1059 *Third Domain*

1060 The third domain is symbolic links, referencing files not of type directory, specified to utilities
1061 that are performing a traversal of a file hierarchy. (This includes symbolic links specified as
1062 command line filename arguments or encountered during the traversal.)

1063 The intention of the Shell and Utilities volume of IEEE Std 1003.1-2001 is that the operation that
1064 the utility is performing is applied to the symbolic link itself, if that operation is applicable to
1065 symbolic links. The reason that the operation is not required is that symbolic links in some
1066 implementations do not have such attributes as a file owner, and therefore the *chown* operation
1067 would be meaningless. If symbolic links on the system have an owner, it is the intention that the
1068 utility *chown* cause the owner of the symbolic link to change. If symbolic links do not have an
1069 owner, the symbolic link should be ignored. Specifically, by default, no change should be made
1070 to the file referenced by the symbolic link.

1071 *Fourth Domain*

1072 The fourth domain is symbolic links referencing files of type directory, specified to utilities that
1073 are performing a traversal of a file hierarchy. (This includes symbolic links specified as
1074 command line filename arguments or encountered during the traversal.)

1075 Most standard utilities do not, by default, indirect into the file hierarchy referenced by the
1076 symbolic link. (The Shell and Utilities volume of IEEE Std 1003.1-2001 uses the informal term
1077 “physical walk” to describe this case. The case where the utility does indirect through the
1078 symbolic link is termed a “logical walk”.)

1079 There are three reasons for the default to be a physical walk:

- 1080 1. With very few exceptions, a physical walk has been the historical default on UNIX systems
1081 supporting symbolic links. Because some utilities (that is, *rm*) must default to a physical
1082 walk, regardless, changing historical practice in this regard would be confusing to users
1083 and needlessly incompatible.
- 1084 2. For systems where symbolic links have the historical file attributes (that is, *owner*, *group*,
1085 *mode*), defaulting to a logical traversal would require the addition of a new option to the
1086 commands to modify the attributes of the link itself. This is painful and more complex
1087 than the alternatives.
- 1088 3. There is a security issue with defaulting to a logical walk. Historically, the command
1089 *chown -R user file* has been safe for the superuser because *setuid* and *setgid* bits were lost
1090 when the ownership of the file was changed. If the walk were logical, changing ownership

1091 would no longer be safe because a user might have inserted a symbolic link pointing to any
1092 file in the tree. Again, this would necessitate the addition of an option to the commands
1093 doing hierarchy traversal to not indirect through the symbolic links, and historical scripts
1094 doing recursive walks would instantly become security problems. While this is mostly an
1095 issue for system administrators, it is preferable to not have different defaults for different
1096 classes of users.

1097 However, the standard developers agreed to leave it unspecified to achieve consensus.

1098 As consistently as possible, users may cause standard utilities performing a file hierarchy
1099 traversal to follow any symbolic links named on the command line, regardless of the type of file
1100 they reference, by specifying the **-H** (for half logical) option. This option is intended to make the
1101 command line name space look like the logical name space.

1102 As consistently as possible, users may cause standard utilities performing a file hierarchy
1103 traversal to follow any symbolic links named on the command line as well as any symbolic links
1104 encountered during the traversal, regardless of the type of file they reference, by specifying the
1105 **-L** (for logical) option. This option is intended to make the entire name space look like the
1106 logical name space.

1107 For consistency, implementors are encouraged to use the **-P** (for “physical”) flag to specify the
1108 physical walk in utilities that do logical walks by default for whatever reason. The only standard
1109 utilities that require the **-P** option are *cd* and *pwd*; see the note below.

1110 When one or more of the **-H**, **-L**, and **-P** flags can be specified, the last one specified determines
1111 the behavior of the utility. This permits users to alias commands so that the default behavior is a
1112 logical walk and then override that behavior on the command line.

1113 *Exceptions in the Third and Fourth Domains*

1114 The *ls* and *rm* utilities are exceptions to these rules. The *rm* utility never follows symbolic links
1115 and does not support the **-H**, **-L**, or **-P** options. Some historical versions of *ls* always followed
1116 symbolic links given on the command line whether the **-L** option was specified or not. Historical
1117 versions of *ls* did not support the **-H** option. In IEEE Std 1003.1-2001, unless one of the **-H** or **-L**
1118 options is specified, the *ls* utility only follows symbolic links to directories that are given as
1119 operands. The *ls* utility does not support the **-P** option.

1120 The Shell and Utilities volume of IEEE Std 1003.1-2001 requires that the standard utilities *ls*, *find*,
1121 and *pax* detect infinite loops when doing logical walks; that is, a directory, or more commonly a
1122 symbolic link, that refers to an ancestor in the current file hierarchy. If the file system itself is
1123 corrupted, causing the infinite loop, it may be impossible to recover. Because *find* and *ls* are often
1124 used in system administration and security applications, they should attempt to recover and
1125 continue as best as they can. The *pax* utility should terminate because the archive it was creating
1126 is by definition corrupted. Other, less vital, utilities should probably simply terminate as well.
1127 Implementations are strongly encouraged to detect infinite loops in all utilities.

1128 Historical practice is shown in Table A-1 (on page 30). The heading **SVID3** stands for the Third
1129 Edition of the System V Interface Definition.

1130 Historically, several shells have had built-in versions of the *pwd* utility. In some of these shells,
1131 *pwd* reported the physical path, and in others, the logical path. Implementations of the shell
1132 corresponding to IEEE Std 1003.1-2001 must report the logical path by default. Earlier versions of
1133 IEEE Std 1003.1-2001 did not require the *pwd* utility to be a built-in utility. Now that *pwd* is
1134 required to set an environment variable in the current shell execution environment, it must be a
1135 built-in utility.

1136 The *cd* command is required, by default, to treat the filename dot-dot logically. Implementors are
1137 required to support the **-P** flag in *cd* so that users can have their current environment handled

1138 physically. In 4.3 BSD, *chgrp* during tree traversal changed the group of the symbolic link, not
 1139 the target. Symbolic links in 4.4 BSD do not have *owner*, *group*, *mode*, or other standard UNIX
 1140 system file attributes.

1141 **Table A-1** Historical Practice for Symbolic Links

Utility	SVID3	4.3 BSD	4.4 BSD	POSIX	Comments
1142 <i>cd</i>				-L	Treat ". ." logically.
1143 <i>cd</i>				-P	Treat ". ." physically.
1144 <i>chgrp</i>			-H	-H	Follow command line symlinks.
1145 <i>chgrp</i>			-h	-L	Follow symlinks.
1146 <i>chgrp</i>	-h			-h	Affect the symlink.
1147 <i>chmod</i>					Affect the symlink.
1148 <i>chmod</i>			-H		Follow command line symlinks.
1149 <i>chmod</i>			-h		Follow symlinks.
1150 <i>chown</i>			-H	-H	Follow command line symlinks.
1151 <i>chown</i>			-h	-L	Follow symlinks.
1152 <i>chown</i>	-h			-h	Affect the symlink.
1153 <i>cp</i>			-H	-H	Follow command line symlinks.
1154 <i>cp</i>			-h	-L	Follow symlinks.
1155 <i>cpio</i>	-L		-L		Follow symlinks.
1156 <i>du</i>			-H	-H	Follow command line symlinks.
1157 <i>du</i>			-h	-L	Follow symlinks.
1158 <i>file</i>	-h			-h	Affect the symlink.
1159 <i>find</i>			-H	-H	Follow command line symlinks.
1160 <i>find</i>			-h	-L	Follow symlinks.
1161 <i>find</i>	-follow		-follow		Follow symlinks.
1162 <i>ln</i>	-s	-s	-s	-s	Create a symbolic link.
1163 <i>ls</i>	-L	-L	-L	-L	Follow symlinks.
1164 <i>ls</i>				-H	Follow command line symlinks.
1165 <i>mv</i>					Operates on the symlink.
1166 <i>pax</i>			-H	-H	Follow command line symlinks.
1167 <i>pax</i>			-h	-L	Follow symlinks.
1168 <i>pwd</i>				-L	Printed path may contain symlinks.
1169 <i>pwd</i>				-P	Printed path will not contain symlinks.
1170 <i>rm</i>					Operates on the symlink.
1171 <i>tar</i>			-H		Follow command line symlinks.
1172 <i>tar</i>		-h	-h		Follow symlinks.
1173 <i>test</i>	-h		-h	-h	Affect the symlink.

1175 Synchronously-Generated Signal

1176 Those signals that may be generated synchronously include SIGABRT, SIGBUS, SIGILL, SIGFPE,
 1177 SIGPIPE, and SIGSEGV.

1178 Any signal sent via the *raise()* function or a *kill()* function targeting the current process is also
 1179 considered synchronous.

1180 System Call*

1181 The distinction between a “system call” and a “library routine” is an implementation detail that
1182 may differ between implementations and has thus been excluded from POSIX.1.

1183 See “Interface, Not Implementation” in **Introduction** (on page xv).

1184 System Reboot

1185 A “system reboot” is an event initiated by an unspecified circumstance that causes all processes
1186 (other than special system processes) to be terminated in an implementation-defined manner,
1187 after which any changes to the state and contents of files created or written to by a Conforming
1188 POSIX.1 Application prior to the event are implementation-defined.

1189 Synchronized I/O Data (and File) Integrity Completion

1190 These terms specify that for synchronized read operations, pending writes must be successfully
1191 completed before the read operation can complete. This is motivated by two circumstances.
1192 Firstly, when synchronizing processes can access the same file, but not share common buffers
1193 (such as for a remote file system), this requirement permits the reading process to guarantee that
1194 it can read data written remotely. Secondly, having data written synchronously is insufficient to
1195 guarantee the order with respect to a subsequent write by a reading process, and thus this extra
1196 read semantic is necessary.

1197 Text File

1198 The term “text file” does not prevent the inclusion of control or other non-printable characters
1199 (other than NUL). Therefore, standard utilities that list text files as inputs or outputs are either
1200 able to process the special characters or they explicitly describe their limitations within their
1201 individual descriptions. The definition of “text file” has caused controversy. The only difference
1202 between text and binary files is that text files have lines of less than {LINE_MAX} bytes, with no
1203 NUL characters, each terminated by a <newline>. The definition allows a file with a single
1204 <newline>, but not a totally empty file, to be called a text file. If a file ends with an incomplete
1205 line it is not strictly a text file by this definition. The <newline> referred to in
1206 IEEE Std 1003.1-2001 is not some generic line separator, but a single character; files created on
1207 systems where they use multiple characters for ends of lines are not portable to all conforming
1208 systems without some translation process unspecified by IEEE Std 1003.1-2001.

1209 Thread

1210 IEEE Std 1003.1-2001 defines a thread to be a flow of control within a process. Each thread has a
1211 minimal amount of private state; most of the state associated with a process is shared among all
1212 of the threads in the process. While most multi-thread extensions to POSIX have taken this
1213 approach, others have made different decisions.

1214 **Note:** The choice to put threads within a process does not constrain implementations to implement
1215 threads in that manner. However, all functions have to behave as though threads share the
1216 indicated state information with the process from which they were created.

1217 Threads need to share resources in order to cooperate. Memory has to be widely shared between
1218 threads in order for the threads to cooperate at a fine level of granularity. Threads keep data
1219 structures and the locks protecting those data structures in shared memory. For a data structure
1220 to be usefully shared between threads, such structures should not refer to any data that can only
1221 be interpreted meaningfully by a single thread. Thus, any system resources that might be
1222 referred to in data structures need to be shared between all threads. File descriptors, pathnames,
1223 and pointers to stack variables are all things that programmers want to share between their
1224 threads. Thus, the file descriptor table, the root directory, the current working directory, and the

- 1225 address space have to be shared.
- 1226 Library implementations are possible as long as the effective behavior is as if system services
1227 invoked by one thread do not suspend other threads. This may be difficult for some library
1228 implementations on systems that do not provide asynchronous facilities.
- 1229 See Section B.2.9 (on page 151) for additional rationale.
- 1230 **Thread ID**
- 1231 See Section B.2.9.2 (on page 167) for additional rationale.
- 1232 **Thread-Safe Function**
- 1233 All functions required by IEEE Std 1003.1-2001 need to be thread-safe; see Section A.4.16 (on
1234 page 40) and Section B.2.9.1 (on page 164) for additional rationale.
- 1235 **User Database**
- 1236 There are no references in IEEE Std 1003.1-2001 to a “passwd file” or a “group file”, and there is
1237 no requirement that the *group* or *passwd* databases be kept in files containing editable text. Many
1238 large timesharing systems use *passwd* databases that are hashed for speed. Certain security
1239 classifications prohibit certain information in the *passwd* database from being publicly readable.
- 1240 The term “encoded” is used instead of “encrypted” in order to avoid the implementation
1241 connotations (such as reversibility or use of a particular algorithm) of the latter term.
- 1242 The *getgrent()*, *setgrent()*, *endgrent()*, *getpwent()*, *setpwent()*, and *endpwent()* functions are not
1243 included as part of the base standard because they provide a linear database search capability
1244 that is not generally useful (the *getpwuid()*, *getpwnam()*, *getgrgid()*, and *getgrnam()* functions are
1245 provided for keyed lookup) and because in certain distributed systems, especially those with
1246 different authentication domains, it may not be possible or desirable to provide an application
1247 with the ability to browse the system databases indiscriminately. They are provided on XSI-
1248 conformant systems due to their historical usage by many existing applications.
- 1249 A change from historical implementations is that the structures used by these functions have
1250 fields of the types **gid_t** and **uid_t**, which are required to be defined in the `<sys/types.h>` header.
1251 IEEE Std 1003.1-2001 requires implementations to ensure that these types are defined by
1252 inclusion of `<grp.h>` and `<pwd.h>`, respectively, without imposing any name space pollution or
1253 errors from redefinition of types.
- 1254 IEEE Std 1003.1-2001 is silent about the content of the strings containing user or group names.
1255 These could be digit strings. IEEE Std 1003.1-2001 is also silent as to whether such digit strings
1256 bear any relationship to the corresponding (numeric) user or group ID.
- 1257 **Database Access**
- 1258 The thread-safe versions of the user and group database access functions return values in user-
1259 supplied buffers instead of possibly using static data areas that may be overwritten by each call.

1260 **Virtual Processor***

1261 The term “virtual processor” was chosen as a neutral term describing all kernel-level
1262 schedulable entities, such as processes, Mach tasks, or lightweight processes. Implementing
1263 threads using multiple processes as virtual processors, or implementing multiplexed threads
1264 above a virtual processor layer, should be possible, provided some mechanism has also been
1265 implemented for sharing state between processes or virtual processors. Many systems may also
1266 wish to provide implementations of threads on systems providing “shared processes” or
1267 “variable-weight processes”. It was felt that exposing such implementation details would
1268 severely limit the type of systems upon which the threads interface could be supported and
1269 prevent certain types of valid implementations. It was also determined that a virtual processor
1270 interface was out of the scope of the Rationale (Informative) volume of IEEE Std 1003.1-2001.

1271 **XSI**

1272 This is introduced to allow IEEE Std 1003.1-2001 to be adopted as an IEEE standard and an Open
1273 Group Technical Standard, serving both the POSIX and the Single UNIX Specification in a core
1274 set of volumes.

1275 The term “XSI” has been used for 10 years in connection with the XPG series and the first and
1276 second versions of the base volumes of the Single UNIX Specification. The XSI margin code was
1277 introduced to denote the extended or more restrictive semantics beyond POSIX that are
1278 applicable to UNIX systems.

1279 **A.4 General Concepts**1280 **A.4.1 Concurrent Execution**

1281 There is no additional rationale provided for this section.

1282 **A.4.2 Directory Protection**

1283 There is no additional rationale provided for this section.

1284 **A.4.3 Extended Security Controls**

1285 Allowing an implementation to define extended security controls enables the use of
1286 IEEE Std 1003.1-2001 in environments that require different or more rigorous security than that
1287 provided in POSIX.1. Extensions are allowed in two areas: privilege and file access permissions.
1288 The semantics of these areas have been defined to permit extensions with reasonable, but not
1289 exact, compatibility with all existing practices. For example, the elimination of the superuser
1290 definition precludes identifying a process as privileged or not by virtue of its effective user ID.

1291 **A.4.4 File Access Permissions**

1292 A process should not try to anticipate the result of an attempt to access data by *a priori* use of
1293 these rules. Rather, it should make the attempt to access data and examine the return value (and
1294 possibly *errno* as well), or use *access()*. An implementation may include other security
1295 mechanisms in addition to those specified in POSIX.1, and an access attempt may fail because of
1296 those additional mechanisms, even though it would succeed according to the rules given in this
1297 section. (For example, the user’s security level might be lower than that of the object of the access
1298 attempt.) The supplementary group IDs provide another reason for a process to not attempt to
1299 anticipate the result of an access attempt.

1300 **A.4.5 File Hierarchy**

1301 Though the file hierarchy is commonly regarded to be a tree, POSIX.1 does not define it as such
1302 for three reasons:

- 1303 1. Links may join branches.
- 1304 2. In some network implementations, there may be no single absolute root directory; see
1305 *pathname resolution*.
- 1306 3. With symbolic links, the file system need not be a tree or even a directed acyclic graph.

1307 **A.4.6 Filenames**

1308 Historically, certain filenames have been reserved. This list includes **core**, **/etc/passwd**, and so
1309 on. Conforming applications should avoid these.

1310 Most historical implementations prohibit case folding in filenames; that is, treating uppercase
1311 and lowercase alphabetic characters as identical. However, some consider case folding desirable:

- 1312 • For user convenience
- 1313 • For ease-of-implementation of the POSIX.1 interface as a hosted system on some popular
1314 operating systems

1315 Variants, such as maintaining case distinctions in filenames, but ignoring them in comparisons,
1316 have been suggested. Methods of allowing escaped characters of the case opposite the default
1317 have been proposed.

1318 Many reasons have been expressed for not allowing case folding, including:

- 1319 • No solid evidence has been produced as to whether case-sensitivity or case-insensitivity is
1320 more convenient for users.
- 1321 • Making case-insensitivity a POSIX.1 implementation option would be worse than either
1322 having it or not having it, because:
- 1323 — More confusion would be caused among users.
- 1324 — Application developers would have to account for both cases in their code.
- 1325 — POSIX.1 implementors would still have other problems with native file systems, such as
1326 short or otherwise constrained filenames or pathnames, and the lack of hierarchical
1327 directory structure.
- 1328 • Case folding is not easily defined in many European languages, both because many of them
1329 use characters outside the US ASCII alphabetic set, and because:
- 1330 — In Spanish, the digraph "ll" is considered to be a single letter, the capitalized form of
1331 which may be either "Ll" or "LL", depending on context.
- 1332 — In French, the capitalized form of a letter with an accent may or may not retain the accent,
1333 depending on the country in which it is written.
- 1334 — In German, the sharp ess may be represented as a single character resembling a Greek
1335 beta (β) in lowercase, but as the digraph "SS" in uppercase.
- 1336 — In Greek, there are several lowercase forms of some letters; the one to use depends on its
1337 position in the word. Arabic has similar rules.
- 1338 • Many East Asian languages, including Japanese, Chinese, and Korean, do not distinguish
1339 case and are sometimes encoded in character sets that use more than one byte per character.

- 1340 • Multiple character codes may be used on the same machine simultaneously. There are
 1341 several ISO character sets for European alphabets. In Japan, several Japanese character codes
 1342 are commonly used together, sometimes even in filenames; this is evidently also the case in
 1343 China. To handle case insensitivity, the kernel would have to at least be able to distinguish
 1344 for which character sets the concept made sense.
- 1345 • The file system implementation historically deals only with bytes, not with characters, except
 1346 for slash and the null byte.
- 1347 • The purpose of POSIX.1 is to standardize the common, existing definition, not to change it.
 1348 Mandating case-insensitivity would make all historical implementations non-standard.
- 1349 • Not only the interface, but also application programs would need to change, counter to the
 1350 purpose of having minimal changes to existing application code.
- 1351 • At least one of the original developers of the UNIX system has expressed objection in the
 1352 strongest terms to either requiring case-insensitivity or making it an option, mostly on the
 1353 basis that POSIX.1 should not hinder portability of application programs across related
 1354 implementations in order to allow compatibility with unrelated operating systems.

1355 Two proposals were entertained regarding case folding in filenames:

- 1356 1. Remove all wording that previously permitted case folding.
- 1357 Rationale Case folding is inconsistent with portable filename character set definition
 1358 and filename definition (all characters except slash and null). No known
 1359 implementations allowing all characters except slash and null also do case
 1360 folding.
- 1361 2. Change “though this practice is not recommended:” to “although this practice is strongly
 1362 discouraged.”
- 1363 Rationale If case folding must be included in POSIX.1, the wording should be stronger
 1364 to discourage the practice.

1365 The consensus selected the first proposal. Otherwise, a conforming application would have to
 1366 assume that case folding would occur when it was not wanted, but that it would not occur when
 1367 it was wanted.

1368 **A.4.7 File Times Update**

1369 This section reflects the actions of historical implementations. The times are not updated
 1370 immediately, but are only marked for update by the functions. An implementation may update
 1371 these times immediately.

1372 The accuracy of the time update values is intentionally left unspecified so that systems can
 1373 control the bandwidth of a possible covert channel.

1374 The wording was carefully chosen to make it clear that there is no requirement that the
 1375 conformance document contain information that might incidentally affect file update times. Any
 1376 function that performs pathname resolution might update several *st_atime* fields. Functions such
 1377 as *getpwnam()* and *getgrnam()* might update the *st_atime* field of some specific file or files. It is
 1378 intended that these are not required to be documented in the conformance document, but they
 1379 should appear in the system documentation.

1380 **A.4.8 Host and Network Byte Order**

1381 There is no additional rationale provided for this section.

1382 **A.4.9 Measurement of Execution Time**

1383 The methods used to measure the execution time of processes and threads, and the precision of
 1384 these measurements, may vary considerably depending on the software architecture of the
 1385 implementation, and on the underlying hardware. Implementations can also make tradeoffs
 1386 between the scheduling overhead and the precision of the execution time measurements.
 1387 IEEE Std 1003.1-2001 does not impose any requirement on the accuracy of the execution time; it
 1388 instead specifies that the measurement mechanism and its precision are implementation-
 1389 defined.

1390 **A.4.10 Memory Synchronization**

1391 In older multi-processors, access to memory by the processors was strictly multiplexed. This
 1392 meant that a processor executing program code interrogates or modifies memory in the order
 1393 specified by the code and that all the memory operation of all the processors in the system
 1394 appear to happen in some global order, though the operation histories of different processors are
 1395 interleaved arbitrarily. The memory operations of such machines are said to be sequentially
 1396 consistent. In this environment, threads can synchronize using ordinary memory operations. For
 1397 example, a producer thread and a consumer thread can synchronize access to a circular data
 1398 buffer as follows:

```

1399     int rdptr = 0;
1400     int wrptr = 0;
1401     data_t buf[BUFSIZE];

1402     Thread 1:
1403         while (work_to_do) {
1404             int next;

1405             buf[wrptr] = produce();
1406             next = (wrptr + 1) % BUFSIZE;
1407             while (rdptr == next)
1408                 ;
1409             wrptr = next;
1410         }

1411     Thread 2:
1412         while (work_to_do) {
1413             while (rdptr == wrptr)
1414                 ;
1415             consume(buf[rdptr]);
1416             rdptr = (rdptr + 1) % BUFSIZE;
1417         }

```

1418 In modern multi-processors, these conditions are relaxed to achieve greater performance. If one
 1419 processor stores values in location A and then location B, then other processors loading data
 1420 from location B and then location A may see the new value of B but the old value of A. The
 1421 memory operations of such machines are said to be weakly ordered. On these machines, the
 1422 circular buffer technique shown in the example will fail because the consumer may see the new
 1423 value of *wrptr* but the old value of the data in the buffer. In such machines, synchronization can
 1424 only be achieved through the use of special instructions that enforce an order on memory
 1425 operations. Most high-level language compilers only generate ordinary memory operations to

1426 take advantage of the increased performance. They usually cannot determine when memory
1427 operation order is important and generate the special ordering instructions. Instead, they rely on
1428 the programmer to use synchronization primitives correctly to ensure that modifications to a
1429 location in memory are ordered with respect to modifications and/or access to the same location
1430 in other threads. Access to read-only data need not be synchronized. The resulting program is
1431 said to be data race-free.

1432 Synchronization is still important even when accessing a single primitive variable (for example,
1433 an integer). On machines where the integer may not be aligned to the bus data width or be larger
1434 than the data width, a single memory load may require multiple memory cycles. This means
1435 that it may be possible for some parts of the integer to have an old value while other parts have a
1436 newer value. On some processor architectures this cannot happen, but portable programs cannot
1437 rely on this.

1438 In summary, a portable multi-threaded program, or a multi-process program that shares
1439 writable memory between processes, has to use the synchronization primitives to synchronize
1440 data access. It cannot rely on modifications to memory being observed by other threads in the
1441 order written in the application or even on modification of a single variable being seen
1442 atomically.

1443 Conforming applications may only use the functions listed to synchronize threads of control
1444 with respect to memory access. There are many other candidates for functions that might also be
1445 used. Examples are: signal sending and reception, or pipe writing and reading. In general, any
1446 function that allows one thread of control to wait for an action caused by another thread of
1447 control is a candidate. IEEE Std 1003.1-2001 does not require these additional functions to
1448 synchronize memory access since this would imply the following:

- 1449 • All these functions would have to be recognized by advanced compilation systems so that
1450 memory operations and calls to these functions are not reordered by optimization.
- 1451 • All these functions would potentially have to have memory synchronization instructions
1452 added, depending on the particular machine.
- 1453 • The additional functions complicate the model of how memory is synchronized and make
1454 automatic data race detection techniques impractical.

1455 Formal definitions of the memory model were rejected as unreadable by the vast majority of
1456 programmers. In addition, most of the formal work in the literature has concentrated on the
1457 memory as provided by the hardware as opposed to the application programmer through the
1458 compiler and runtime system. It was believed that a simple statement intuitive to most
1459 programmers would be most effective. IEEE Std 1003.1-2001 defines functions that can be used
1460 to synchronize access to memory, but it leaves open exactly how one relates those functions to
1461 the semantics of each function as specified elsewhere in IEEE Std 1003.1-2001.
1462 IEEE Std 1003.1-2001 also does not make a formal specification of the partial ordering in time
1463 that the functions can impose, as that is implied in the description of the semantics of each
1464 function. It simply states that the programmer has to ensure that modifications do not occur
1465 “simultaneously” with other access to a memory location.

1466 **A.4.11 Pathname Resolution**

1467 It is necessary to differentiate between the definition of pathname and the concept of pathname
 1468 resolution with respect to the handling of trailing slashes. By specifying the behavior here, it is
 1469 not possible to provide an implementation that is conforming but extends all interfaces that
 1470 handle pathnames to also handle strings that are not legal pathnames (because they have trailing
 1471 slashes).

1472 Pathnames that end with one or more trailing slash characters must refer to directory paths.
 1473 Previous versions of IEEE Std 1003.1-2001 were not specific about the distinction between
 1474 trailing slashes on files and directories, and both were permitted.

1475 Two types of implementation have been prevalent; those that ignored trailing slash characters
 1476 on all pathnames regardless, and those that permitted them only on existing directories.

1477 IEEE Std 1003.1-2001 requires that a pathname with a trailing slash character be treated as if it
 1478 had a trailing "/" everywhere.

1479 Note that this change does not break any conforming applications; since there were two
 1480 different types of implementation, no application could have portably depended on either
 1481 behavior. This change does however require some implementations to be altered to remain
 1482 compliant. Substantial discussion over a three-year period has shown that the benefits to
 1483 application developers outweighs the disadvantages for some vendors.

1484 On a historical note, some early applications automatically appended a "/" to every path.
 1485 Rather than fix the applications, the system implementation was modified to accept this
 1486 behavior by ignoring any trailing slash.

1487 Each directory has exactly one parent directory which is represented by the name **dot-dot** in the
 1488 first directory. No other directory, regardless of linkages established by symbolic links, is
 1489 considered the parent directory by IEEE Std 1003.1-2001.

1490 There are two general categories of interfaces involving pathname resolution: those that follow
 1491 the symbolic link, and those that do not. There are several exceptions to this rule; for example,
 1492 *open(path, O_CREAT|O_EXCL)* will fail when *path* names a symbolic link. However, in all other
 1493 situations, the *open()* function will follow the link.

1494 What the filename **dot-dot** refers to relative to the root directory is implementation-defined. In
 1495 Version 7 it refers to the root directory itself; this is the behavior mentioned in
 1496 IEEE Std 1003.1-2001. In some networked systems the construction *./hostname/* is used to refer
 1497 to the root directory of another host, and POSIX.1 permits this behavior.

1498 Other networked systems use the construct *//hostname* for the same purpose; that is, a double
 1499 initial slash is used. There is a potential problem with existing applications that create full
 1500 pathnames by taking a trunk and a relative pathname and making them into a single string
 1501 separated by '/', because they can accidentally create networked pathnames when the trunk is
 1502 '/'. This practice is not prohibited because such applications can be made to conform by
 1503 simply changing to use "/" as a separator instead of '/':

- 1504 • If the trunk is '/', the full pathname will begin with "///" (the initial '/' and the
 1505 separator '/'). This is the same as '/', which is what is desired. (This is the general case
 1506 of making a relative pathname into an absolute one by prefixing with "///" instead of '/'.)
- 1507 • If the trunk is "/A", the result is "/A//..."; since non-leading sequences of two or more
 1508 slashes are treated as a single slash, this is equivalent to the desired "/A//...".
- 1509 • If the trunk is "//A", the implementation-defined semantics will apply. (The multiple slash
 1510 rule would apply.)

1511 Application developers should avoid generating pathnames that start with `"/"`.
1512 Implementations are strongly encouraged to avoid using this special interpretation since a
1513 number of applications currently do not follow this practice and may inadvertently generate
1514 `"/. . ."`.

1515 The term “root directory” is only defined in POSIX.1 relative to the process. In some
1516 implementations, there may be no absolute root directory. The initialization of the root directory
1517 of a process is implementation-defined.

1518 **A.4.12 Process ID Reuse**

1519 There is no additional rationale provided for this section.

1520 **A.4.13 Scheduling Policy**

1521 There is no additional rationale provided for this section.

1522 **A.4.14 Seconds Since the Epoch**

1523 Coordinated Universal Time (UTC) includes leap seconds. However, in POSIX time (seconds
1524 since the Epoch), leap seconds are ignored (not applied) to provide an easy and compatible
1525 method of computing time differences. Broken-down POSIX time is therefore not necessarily
1526 UTC, despite its appearance.

1527 As of September 2000, 24 leap seconds had been added to UTC since the Epoch, 1 January, 1970.
1528 Historically, one leap second is added every 15 months on average, so this offset can be expected
1529 to grow steadily with time.

1530 Most systems’ notion of “time” is that of a continuously increasing value, so this value should
1531 increase even during leap seconds. However, not only do most systems not keep track of leap
1532 seconds, but most systems are probably not synchronized to any standard time reference.
1533 Therefore, it is inappropriate to require that a time represented as seconds since the Epoch
1534 precisely represent the number of seconds between the referenced time and the Epoch.

1535 It is sufficient to require that applications be allowed to treat this time as if it represented the
1536 number of seconds between the referenced time and the Epoch. It is the responsibility of the
1537 vendor of the system, and the administrator of the system, to ensure that this value represents
1538 the number of seconds between the referenced time and the Epoch as closely as necessary for the
1539 application being run on that system.

1540 It is important that the interpretation of time names and seconds since the Epoch values be
1541 consistent across conforming systems; that is, it is important that all conforming systems
1542 interpret “536 457 599 seconds since the Epoch” as 59 seconds, 59 minutes, 23 hours 31 December
1543 1986, regardless of the accuracy of the system’s idea of the current time. The expression is given
1544 to ensure a consistent interpretation, not to attempt to specify the calendar. The relationship
1545 between `tm_yday` and the day of week, day of month, and month is in accordance with the
1546 Gregorian calendar, and so is not specified in POSIX.1.

1547 Consistent interpretation of seconds since the Epoch can be critical to certain types of distributed
1548 applications that rely on such timestamps to synchronize events. The accrual of leap seconds in
1549 a time standard is not predictable. The number of leap seconds since the Epoch will likely
1550 increase. POSIX.1 is more concerned about the synchronization of time between applications of
1551 astronomically short duration.

1552 Note that `tm_yday` is zero-based, not one-based, so the day number in the example above is 364.
1553 Note also that the division is an integer division (discarding remainder) as in the C language.

1554 Note also that the meaning of *gmtime()*, *localtime()*, and *mktime()* is specified in terms of this
1555 expression. However, the ISO C standard computes *tm_yday* from *tm_mday*, *tm_mon*, and
1556 *tm_year* in *mktime()*. Because it is stated as a (bidirectional) relationship, not a function, and
1557 because the conversion between month-day-year and day-of-year dates is presumed well known
1558 and is also a relationship, this is not a problem.

1559 Implementations that implement **time_t** as a signed 32-bit integer will overflow in 2038. The
1560 data size for **time_t** is as per the ISO C standard definition, which is implementation-defined.

1561 See also **Epoch** (on page 17).

1562 The topic of whether seconds since the Epoch should account for leap seconds has been debated
1563 on a number of occasions, and each time consensus was reached (with acknowledged dissent
1564 each time) that the majority of users are best served by treating all days identically. (That is, the
1565 majority of applications were judged to assume a single length—as measured in seconds since
1566 the Epoch—for all days. Thus, leap seconds are not applied to seconds since the Epoch.) Those
1567 applications which do care about leap seconds can determine how to handle them in whatever
1568 way those applications feel is best. This was particularly emphasized because there was
1569 disagreement about what the best way of handling leap seconds might be. It is a practical
1570 impossibility to mandate that a conforming implementation must have a fixed relationship to
1571 any particular official clock (consider isolated systems, or systems performing “reruns” by
1572 setting the clock to some arbitrary time).

1573 Note that as a practical consequence of this, the length of a second as measured by some external
1574 standard is not specified. This unspecified second is nominally equal to an International System
1575 (SI) second in duration. Applications must be matched to a system that provides the particular
1576 handling of external time in the way required by the application.

1577 **A.4.15 Semaphore**

1578 There is no additional rationale provided for this section.

1579 **A.4.16 Thread-Safety**

1580 Where the interface of a function required by IEEE Std 1003.1-2001 precludes thread-safety, an
1581 alternate thread-safe form is provided. The names of these thread-safe forms are the same as the
1582 non-thread-safe forms with the addition of the suffix “_r”. The suffix “_r” is historical, where
1583 the ‘r’ stood for “reentrant”.

1584 In some cases, thread-safety is provided by restricting the arguments to an existing function.

1585 See also Section B.2.9.1 (on page 164).

1586 **A.4.17 Tracing**

1587 Refer to Section B.2.11 (on page 180).

1588 **A.4.18 Treatment of Error Conditions for Mathematical Functions**

1589 There is no additional rationale provided for this section.

1590 **A.4.19 Treatment of NaN Arguments for Mathematical Functions**

1591 There is no additional rationale provided for this section.

1592 **A.4.20 Utility**

1593 There is no additional rationale provided for this section.

1594 **A.4.21 Variable Assignment**

1595 There is no additional rationale provided for this section.

1596 **A.5 File Format Notation**

1597 The notation for spaces allows some flexibility for application output. Note that an empty
 1598 character position in *format* represents one or more <blank>s on the output (not *white space*,
 1599 which can include <newline>s). Therefore, another utility that reads that output as its input
 1600 must be prepared to parse the data using *scanf()*, *awk*, and so on. The 'Δ' character is used when
 1601 exactly one <space> is output.

1602 The treatment of integers and spaces is different from the *printf()* function in that they can be
 1603 surrounded with <blank>s. This was done so that, given a format such as:

1604 `"%d\n", <foo>`1605 the implementation could use a *printf()* call such as:1606 `printf("%6d\n", foo);`1607 and still conform. This notation is thus somewhat like *scanf()* in addition to *printf()*.

1608 The *printf()* function was chosen as a model because most of the standard developers were
 1609 familiar with it. One difference from the C function *printf()* is that the l and h conversion
 1610 specifier characters are not used. As expressed by the Shell and Utilities volume of
 1611 IEEE Std 1003.1-2001, there is no differentiation between decimal values for type **int**, type **long**,
 1612 or type **short**. The conversion specifications %d or %i should be interpreted as an arbitrary
 1613 length sequence of digits. Also, no distinction is made between single precision and double
 1614 precision numbers (**float** or **double** in C). These are simply referred to as floating-point numbers.

1615 Many of the output descriptions in the Shell and Utilities volume of IEEE Std 1003.1-2001 use the
 1616 term "line", such as:

1617 `"%s", <input line>`

1618 Since the definition of *line* includes the trailing <newline> already, there is no need to include a
 1619 '\n' in the format; a double <newline> would otherwise result.

1620 **A.6 Character Set**

1621 **A.6.1 Portable Character Set**

1622 The portable character set is listed in full so there is no dependency on the ISO/IEC 646:1991
1623 standard (or historically ASCII) encoded character set, although the set is identical to the
1624 characters defined in the International Reference version of the ISO/IEC 646:1991 standard.

1625 IEEE Std 1003.1-2001 poses no requirement that multiple character sets or codesets be supported,
1626 leaving this as a marketing differentiation for implementors. Although multiple charmap files
1627 are supported, it is the responsibility of the implementation to provide the file(s); if only one is
1628 provided, only that one will be accessible using the *localedef -f* option.

1629 The statement about invariance in codesets for the portable character set is worded to avoid
1630 precluding implementations where multiple incompatible codesets are available (for instance,
1631 ASCII and EBCDIC). The standard utilities cannot be expected to produce predictable results if
1632 they access portable characters that vary on the same implementation.

1633 Not all character sets need include the portable character set, but each locale must include it. For
1634 example, a Japanese-based locale might be supported by a mixture of character sets: JIS X 0201
1635 Roman (a Japanese version of the ISO/IEC 646:1991 standard), JIS X 0208, and JIS X 0201
1636 Katakana. Not all of these character sets include the portable characters, but at least one does
1637 (JIS X 0201 Roman).

1638 **A.6.2 Character Encoding**

1639 Encoding mechanisms based on single shifts, such as the EUC encoding used in some Asian and
1640 other countries, can be supported via the current charmap mechanism. With single-shift
1641 encoding, each character is preceded by a shift code (SS2 or SS3). A complete EUC code,
1642 consisting of the portable character set (G0) and up to three additional character sets (G1, G2,
1643 G3), can be described using the current charmap mechanism; the encoding for each character in
1644 additional character sets G2 and G3 must then include their single-shift code. Other mechanisms
1645 to support locales based on encoding mechanisms such as locking shift are not addressed by this
1646 volume of IEEE Std 1003.1-2001.

1647 **A.6.3 C Language Wide-Character Codes**

1648 There is no additional rationale provided for this section.

1649 **A.6.4 Character Set Description File**

1650 IEEE PASC Interpretation 1003.2 #196 is applied, removing three lines of text dealing with
1651 ranges of symbolic names using position constant values which had been erroneously included
1652 in the final IEEE P1003.2b draft standard.

1653 *A.6.4.1 State-Dependent Character Encodings*

1654 A requirement was considered that would force utilities to eliminate any redundant locking
1655 shifts, but this was left as a quality of implementation issue.

1656 This change satisfies the following requirement from the ISO POSIX-2:1993 standard, Annex
1657 H.1:

1658 The support of state-dependent (shift encoding) character sets should be addressed fully. See
 1659 descriptions of these in the Base Definitions volume of IEEE Std 1003.1-2001, Section 6.2, Character
 1660 Encoding. If such character encodings are supported, it is expected that this will impact the Base
 1661 Definitions volume of IEEE Std 1003.1-2001, Section 6.2, Character Encoding, the Base Definitions
 1662 volume of IEEE Std 1003.1-2001, Chapter 7, Locale, the Base Definitions volume of
 1663 IEEE Std 1003.1-2001, Chapter 9, Regular Expressions, and the comm, cut, diff, grep, head, join,
 1664 paste, and tail utilities.

1665 The character set description file provides:

- 1666 • The capability to describe character set attributes (such as collation order or character
 1667 classes) independent of character set encoding, and using only the characters in the portable
 1668 character set. This makes it possible to create generic *localedef* source files for all codesets that
 1669 share the portable character set (such as the ISO 8859 family or IBM Extended ASCII).
- 1670 • Standardized symbolic names for all characters in the portable character set, making it
 1671 possible to refer to any such character regardless of encoding.

1672 Implementations are free to choose their own symbolic names, as long as the names identified
 1673 by the Base Definitions volume of IEEE Std 1003.1-2001 are also defined; this provides support
 1674 for already existing “character names”.

1675 The names selected for the members of the portable character set follow the
 1676 ISO/IEC 8859-1:1998 standard and the ISO/IEC 10646-1:2000 standard. However, several
 1677 commonly used UNIX system names occur as synonyms in the list:

- 1678 • The historical UNIX system names are used for control characters.
- 1679 • The word “slash” is given in addition to “solidus”.
- 1680 • The word “backslash” is given in addition to “reverse-solidus”.
- 1681 • The word “hyphen” is given in addition to “hyphen-minus”.
- 1682 • The word “period” is given in addition to “full-stop”.
- 1683 • For digits, the word “digit” is eliminated.
- 1684 • For letters, the words “Latin Capital Letter” and “Latin Small Letter” are eliminated.
- 1685 • The words “left brace” and “right brace” are given in addition to “left-curly-bracket” and
 1686 “right-curly-bracket”.
- 1687 • The names of the digits are preferred over the numbers to avoid possible confusion between
 1688 ‘0’ and ‘o’, and between ‘1’ and ‘l’ (one and the letter ell).

1689 The names for the control characters in the Base Definitions volume of IEEE Std 1003.1-2001,
 1690 Chapter 6, Character Set were taken from the ISO/IEC 4873:1991 standard.

1691 The charmap file was introduced to resolve problems with the portability of, especially, *localedef*
 1692 sources. IEEE Std 1003.1-2001 assumes that the portable character set is constant across all
 1693 locales, but does not prohibit implementations from supporting two incompatible codings, such
 1694 as both ASCII and EBCDIC. Such dual-support implementations should have all charmaps and
 1695 *localedef* sources encoded using one portable character set, in effect cross-compiling for the other
 1696 environment. Naturally, charmaps (and *localedef* sources) are only portable without
 1697 transformation between systems using the same encodings for the portable character set. They
 1698 can, however, be transformed between two sets using only a subset of the actual characters (the
 1699 portable character set). However, the particular coded character set used for an application or an
 1700 implementation does not necessarily imply different characteristics or collation; on the contrary,
 1701 these attributes should in many cases be identical, regardless of codeset. The charmap provides

1702 the capability to define a common locale definition for multiple codesets (the same *localedef*
 1703 source can be used for codesets with different extended characters; the ability in the charmap to
 1704 define empty names allows for characters missing in certain codesets).

1705 The `<escape_char>` declaration was added at the request of the international community to ease
 1706 the creation of portable charmap files on terminals not implementing the default backslash
 1707 escape. The `<comment_char>` declaration was added at the request of the international
 1708 community to eliminate the potential confusion between the number sign and the pound sign.

1709 The octal number notation with no leading zero required was selected to match those of *awk* and
 1710 *tr* and is consistent with that used by *localedef*. To avoid confusion between an octal constant
 1711 and the back-references used in *localedef* source, the octal, hexadecimal, and decimal constants
 1712 must contain at least two digits. As single-digit constants are relatively rare, this should not
 1713 impose any significant hardship. Provision is made for more digits to account for systems in
 1714 which the byte size is larger than 8 bits. For example, a Unicode (ISO/IEC 10646-1:2000
 1715 standard) system that has defined 16-bit bytes may require six octal, four hexadecimal, and five
 1716 decimal digits.

1717 The decimal notation is supported because some newer international standards define character
 1718 values in decimal, rather than in the old column/row notation.

1719 The charmap identifies the coded character sets supported by an implementation. At least one
 1720 charmap must be provided, but no implementation is required to provide more than one.
 1721 Likewise, implementations can allow users to generate new charmaps (for instance, for a new
 1722 version of the ISO 8859 family of coded character sets), but does not have to do so. If users are
 1723 allowed to create new charmaps, the system documentation describes the rules that apply (for
 1724 instance, “only coded character sets that are supersets of the ISO/IEC 646: 1991 standard IRV, no
 1725 multi-byte characters”).

1726 This addition of the **WIDTH** specification satisfies the following requirement from the
 1727 ISO POSIX-2: 1993 standard, Annex H.1:

1728 (9) *The definition of column position relies on the implementation’s knowledge of the integral width*
 1729 *of the characters. The charmap or LC_CTYPE locale definitions should be enhanced to allow*
 1730 *application specification of these widths.*

1731 The character “width” information was first considered for inclusion under *LC_CTYPE* but was
 1732 moved because it is more closely associated with the information in the charmap than
 1733 information in the locale source (cultural conventions information). Concerns were raised that
 1734 formalizing this type of information is moving the locale source definition from the codeset-
 1735 independent entity that it was designed to be to a repository of codeset-specific information. A
 1736 similar issue occurred with the `<code_set_name>`, `<mb_cur_max>`, and `<mb_cur_min>`
 1737 information, which was resolved to reside in the charmap definition.

1738 The width definition was added to the IEEE P1003.2b draft standard with the intent that the
 1739 *wcswidth()* and/or *wcwidth()* functions (currently specified in the System Interfaces volume of
 1740 IEEE Std 1003.1-2001) be the mechanism to retrieve the character width information.

1741 A.7 Locale

1742 A.7.1 General

1743 The description of locales is based on work performed in the UniForum Technical Committee,
1744 Subcommittee on Internationalization. Wherever appropriate, keywords are taken from the
1745 ISO C standard or the X/Open Portability Guide.

1746 The value used to specify a locale with environment variables is the name specified as the *name*
1747 operand to the *localedef* utility when the locale was created. This provides a verifiable method to
1748 create and invoke a locale.

1749 The “object” definitions need not be portable, as long as “source” definitions are. Strictly
1750 speaking, source definitions are portable only between implementations using the same
1751 character set(s). Such source definitions, if they use symbolic names only, easily can be ported
1752 between systems using different codesets, as long as the characters in the portable character set
1753 (see the Base Definitions volume of IEEE Std 1003.1-2001, Section 6.1, Portable Character Set)
1754 have common values between the codesets; this is frequently the case in historical
1755 implementations. Of source, this requires that the symbolic names used for characters outside
1756 the portable character set be identical between character sets. The definition of symbolic names
1757 for characters is outside the scope of IEEE Std 1003.1-2001, but is certainly within the scope of
1758 other standards organizations.

1759 Applications can select the desired locale by invoking the *setlocale()* function (or equivalent)
1760 with the appropriate value. If the function is invoked with an empty string, the value of the
1761 corresponding environment variable is used. If the environment variable is not set or is set to the
1762 empty string, the implementation sets the appropriate environment as defined in the Base
1763 Definitions volume of IEEE Std 1003.1-2001, Chapter 8, Environment Variables.

1764 A.7.2 POSIX Locale

1765 The POSIX locale is equal to the C locale. To avoid being classified as a C-language function, the
1766 name has been changed to the POSIX locale; the environment variable value can be either
1767 "POSIX" or, for historical reasons, "C".

1768 The POSIX definitions mirror the historical UNIX system behavior.

1769 The use of symbolic names for characters in the tables does not imply that the POSIX locale must
1770 be described using symbolic character names, but merely that it may be advantageous to do so.

1771 A.7.3 Locale Definition

1772 The decision to separate the file format from the *localedef* utility description was only partially
1773 editorial. Implementations may provide other interfaces than *localedef*. Requirements on “the
1774 utility”, mostly concerning error messages, are described in this way because they are meant to
1775 affect the other interfaces implementations may provide as well as *localedef*.

1776 The text about POSIX2_LOCALEDEF does not mean that internationalization is optional; only
1777 that the functionality of the *localedef* utility is. REs, for instance, must still be able to recognize,
1778 for example, character class expressions such as "[[:alpha:]]". A possible analogy is with
1779 an applications development environment; while all conforming implementations must be
1780 capable of executing applications, not all need to have the development environment installed.
1781 The assumption is that the capability to modify the behavior of utilities (and applications) via
1782 locale settings must be supported. If the *localedef* utility is not present, then the only choice is to
1783 select an existing (presumably implementation-documented) locale. An implementation could,
1784 for example, choose to support only the POSIX locale, which would in effect limit the amount of

1785 changes from historical implementations quite drastically. The *localedef* utility is still required,
 1786 but would always terminate with an exit code indicating that no locale could be created.
 1787 Supported locales must be documented using the syntax defined in this chapter. (This ensures
 1788 that users can accurately determine what capabilities are provided. If the implementation
 1789 decides to provide additional capabilities to the ones in this chapter, that is already provided
 1790 for.)

1791 If the option is present (that is, locales can be created), then the *localedef* utility must be capable
 1792 of creating locales based on the syntax and rules defined in this chapter. This does not mean that
 1793 the implementation cannot also provide alternate means for creating locales.

1794 The octal, decimal, and hexadecimal notations are the same employed by the charmap facility
 1795 (see the Base Definitions volume of IEEE Std 1003.1-2001, Section 6.4, Character Set Description
 1796 File). To avoid confusion between an octal constant and a back-reference, the octal, hexadecimal,
 1797 and decimal constants must contain at least two digits. As single-digit constants are relatively
 1798 rare, this should not impose any significant hardship. Provision is made for more digits to
 1799 account for systems in which the byte size is larger than 8 bits. For example, a Unicode (see the
 1800 ISO/IEC 10646-1:2000 standard) system that has defined 16-bit bytes may require six octal, four
 1801 hexadecimal, and five decimal digits. As with the charmap file, multi-byte characters are
 1802 described in the locale definition file using “big-endian” notation for reasons of portability.
 1803 There is no requirement that the internal representation in the computer memory be in this same
 1804 order.

1805 One of the guidelines used for the development of this volume of IEEE Std 1003.1-2001 is that
 1806 characters outside the invariant part of the ISO/IEC 646:1991 standard should not be used in
 1807 portable specifications. The backslash character is not in the invariant part; the number sign is,
 1808 but with multiple representations: as a number sign, and as a pound sign. As far as general
 1809 usage of these symbols, they are covered by the “grandfather clause”, but for newly defined
 1810 interfaces, the WG15 POSIX working group has requested that POSIX provide alternate
 1811 representations. Consequently, while the default escape character remains the backslash and the
 1812 default comment character is the number sign, implementations are required to recognize
 1813 alternative representations, identified in the applicable source file via the `<escape_char>` and
 1814 `<comment_char>` keywords.

1815 A.7.3.1 *LC_CTYPE*

1816 The *LC_CTYPE* category is primarily used to define the encoding-independent aspects of a
 1817 character set, such as character classification. In addition, certain encoding-dependent
 1818 characteristics are also defined for an application via the *LC_CTYPE* category.
 1819 IEEE Std 1003.1-2001 does not mandate that the encoding used in the locale is the same as the
 1820 one used by the application because an implementation may decide that it is advantageous to
 1821 define locales in a system-wide encoding rather than having multiple, logically identical locales
 1822 in different encodings, and to convert from the application encoding to the system-wide
 1823 encoding on usage. Other implementations could require encoding-dependent locales.

1824 In either case, the *LC_CTYPE* attributes that are directly dependent on the encoding, such as
 1825 `<mb_cur_max>` and the display width of characters, are not user-specifiable in a locale source
 1826 and are consequently not defined as keywords.

1827 Implementations may define additional keywords or extend the *LC_CTYPE* mechanism to allow
 1828 application-defined keywords.

1829 The text “The ellipsis specification shall only be valid within a single encoded character set” is
 1830 present because it is possible to have a locale supported by multiple character encodings, as
 1831 explained in the rationale for the Base Definitions volume of IEEE Std 1003.1-2001, Section 6.1,
 1832 Portable Character Set. An example given there is of a possible Japanese-based locale supported

1833 by a mixture of the character sets JIS X 0201 Roman, JIS X 0208, and JIS X 0201 Katakana.
 1834 Attempting to express a range of characters across these sets is not logical and the
 1835 implementation is free to reject such attempts.

1836 As the `LC_CTYPE` character classes are based on the ISO C standard character class definition,
 1837 the category does not support multi-character elements. For instance, the German character
 1838 <sharp-s> is traditionally classified as a lowercase letter. There is no corresponding uppercase
 1839 letter; in proper capitalization of German text, the <sharp-s> will be replaced by "SS"; that is, by
 1840 two characters. This kind of conversion is outside the scope of the **toupper** and **tolower**
 1841 keywords.

1842 Where IEEE Std 1003.1-2001 specifies that only certain characters can be specified, as for the
 1843 keywords **digit** and **xdigit**, the specified characters must be from the portable character set, as
 1844 shown. As an example, only the Arabic digits 0 through 9 are acceptable as digits.

1845 The character classes **digit**, **xdigit**, **lower**, **upper**, and **space** have a set of automatically included
 1846 characters. These only need to be specified if the character values (that is, encoding) differs from
 1847 the implementation default values. It is not possible to define a locale without these
 1848 automatically included characters unless some implementation extension is used to prevent
 1849 their inclusion. Such a definition would not be a proper superset of the C locale, and thus, it
 1850 might not be possible for the standard utilities to be implemented as programs conforming to
 1851 the ISO C standard.

1852 The definition of character class **digit** requires that only ten characters—the ones defining
 1853 digits—can be specified; alternate digits (for example, Hindi or Kanji) cannot be specified here.
 1854 However, the encoding may vary if an implementation supports more than one encoding.

1855 The definition of character class **xdigit** requires that the characters included in character class
 1856 **digit** are included here also and allows for different symbols for the hexadecimal digits 10
 1857 through 15.

1858 The inclusion of the **charclass** keyword satisfies the following requirement from the
 1859 ISO POSIX-2: 1993 standard, Annex H.1:

1860 (3) *The `LC_CTYPE` (2.5.2.1) locale definition should be enhanced to allow user-specified additional*
 1861 *character classes, similar in concept to the ISO C standard Multibyte Support Extension (MSE)*
 1862 *iswctype() function.*

1863 This keyword was previously included in The Open Group specifications and is now mandated
 1864 in the Shell and Utilities volume of IEEE Std 1003.1-2001.

1865 The symbolic constant `{CHARCLASS_NAME_MAX}` was also adopted from The Open Group
 1866 specifications. Applications portability is enhanced by the use of symbolic constants.

1867 A.7.3.2 `LC_COLLATE`

1868 The rules governing collation depend to some extent on the use. At least five different levels of
 1869 increasingly complex collation rules can be distinguished:

- 1870 1. *Byte/machine code order*: This is the historical collation order in the UNIX system and many
 1871 proprietary operating systems. Collation is here performed character by character, without
 1872 any regard to context. The primary virtue is that it usually is quite fast and also
 1873 completely deterministic; it works well when the native machine collation sequence
 1874 matches the user expectations.
- 1875 2. *Character order*: On this level, collation is also performed character by character, without
 1876 regard to context. The order between characters is, however, not determined by the code
 1877 values, but on the expectations by the user of the “correct” order between characters. In

1878 addition, such a (simple) collation order can specify that certain characters collate equally
1879 (for example, uppercase and lowercase letters).

1880 3. *String ordering*: On this level, entire strings are compared based on relatively
1881 straightforward rules. Several “passes” may be required to determine the order between
1882 two strings. Characters may be ignored in some passes, but not in others; the strings may
1883 be compared in different directions; and simple string substitutions may be performed
1884 before strings are compared. This level is best described as “dictionary” ordering; it is
1885 based on the spelling, not the pronunciation, or meaning, of the words.

1886 4. *Text search ordering*: This is a further refinement of the previous level, best described as
1887 “telephone book ordering”; some common homonyms (words spelled differently but with
1888 the same pronunciation) are collated together; numbers are collated as if they were spelled
1889 out, and so on.

1890 5. *Semantic-level ordering*: Words and strings are collated based on their meaning; entire words
1891 (such as “the”) are eliminated; the ordering is not deterministic. This usually requires
1892 special software and is highly dependent on the intended use.

1893 While the historical collation order formally is at level 1, for the English language it corresponds
1894 roughly to elements at level 2. The user expects to see the output from the *ls* utility sorted very
1895 much as it would be in a dictionary. While telephone book ordering would be an optimal goal
1896 for standard collation, this was ruled out as the order would be language-dependent.
1897 Furthermore, a requirement was that the order must be determined solely from the text string
1898 and the collation rules; no external information (for example, “pronunciation dictionaries”)
1899 could be required.

1900 As a result, the goal for the collation support is at level 3. This also matches the requirements for
1901 the Canadian collation order, as well as other, known collation requirements for alphabetic
1902 scripts. It specifically rules out collation based on pronunciation rules or based on semantic
1903 analysis of the text.

1904 The syntax for the *LC_COLLATE* category source meets the requirements for level 3 and has
1905 been verified to produce the correct result with examples based on French, Canadian, and
1906 Danish collation order. Because it supports multi-character collating elements, it is also capable
1907 of supporting collation in codesets where a character is expressed using non-spacing characters
1908 followed by the base character (such as the ISO/IEC 6937: 1994 standard).

1909 The directives that can be specified in an operand to the **order_start** keyword are based on the
1910 requirements specified in several proposed standards and in customary use. The following is a
1911 rephrasing of rules defined for “lexical ordering in English and French” by the Canadian
1912 Standards Association (the text in square brackets is rephrased):

- 1913 • Once special characters [punctuation] have been removed from original strings, the ordering
1914 is determined by scanning forwards (left to right) [disregarding case and diacriticals].
- 1915 • In case of equivalence, special characters are once again removed from original strings and
1916 the ordering is determined by scanning backwards (starting from the rightmost character of
1917 the string and back), character by character [disregarding case but considering diacriticals].
- 1918 • In case of repeated equivalence, special characters are removed again from original strings
1919 and the ordering is determined by scanning forwards, character by character [considering
1920 both case and diacriticals].
- 1921 • If there is still an ordering equivalence after the first three rules have been applied, then only
1922 special characters and the position they occupy in the string are considered to determine
1923 ordering. The string that has a special character in the lowest position comes first. If two
1924 strings have a special character in the same position, the character [with the lowest collation

- 1925 value] comes first. In case of equality, the other special characters are considered until there
1926 is a difference or until all special characters have been exhausted.
- 1927 It is estimated that this part of IEEE Std 1003.1-2001 covers the requirements for all European
1928 languages, and no particular problems are anticipated with Slavic or Middle East character sets.
- 1929 The Far East (particularly Japanese/Chinese) collations are often based on contextual
1930 information and pronunciation rules (the same ideogram can have different meanings and
1931 different pronunciations). Such collation, in general, falls outside the desired goal of
1932 IEEE Std 1003.1-2001. There are, however, several other collation rules (stroke/radical or “most
1933 common pronunciation”) that can be supported with the mechanism described here.
- 1934 The character order is defined by the order in which characters and elements are specified
1935 between the **order_start** and **order_end** keywords. Weights assigned to the characters and
1936 elements define the collation sequence; in the absence of weights, the character order is also the
1937 collation sequence.
- 1938 The **position** keyword provides the capability to consider, in a compare, the relative position of
1939 characters not subject to **IGNORE**. As an example, consider the two strings "o-ring" and
1940 "or-ing". Assuming the hyphen is subject to **IGNORE** on the first pass, the two strings
1941 compare equal, and the position of the hyphen is immaterial. On second pass, all characters
1942 except the hyphen are subject to **IGNORE**, and in the normal case the two strings would again
1943 compare equal. By taking position into account, the first collates before the second.
- 1944 **A.7.3.3 LC_MONETARY**
- 1945 The currency symbol does not appear in *LC_MONETARY* because it is not defined in the C locale
1946 of the ISO C standard.
- 1947 The ISO C standard limits the size of decimal points and thousands delimiters to single-byte
1948 values. In locales based on multi-byte coded character sets, this cannot be enforced;
1949 IEEE Std 1003.1-2001 does not prohibit such characters, but makes the behavior unspecified (in
1950 the text “In contexts where other standards ...”).
- 1951 The grouping specification is based on, but not identical to, the ISO C standard. The -1 indicates
1952 that no further grouping is performed; the equivalent of {CHAR_MAX} in the ISO C standard.
- 1953 The text “the value is not available in the locale” is taken from the ISO C standard and is used
1954 instead of the “unspecified” text in early proposals. There is no implication that omitting these
1955 keywords or assigning them values of " " or -1 produces unspecified results; such omissions or
1956 assignments eliminate the effects described for the keyword or produce zero-length strings, as
1957 appropriate.
- 1958 The locale definition is an extension of the ISO C standard *localeconv()* specification. In
1959 particular, rules on how **currency_symbol** is treated are extended to also cover **int_curr_symbol**,
1960 and **p_set_by_space** and **n_sep_by_space** have been augmented with the value 2, which places
1961 a <space> between the sign and the symbol (if they are adjacent; otherwise, it should be treated
1962 as a 0). The following table shows the result of various combinations:

1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975

		p_sep_by_space		
		2	1	0
p_cs_precedes = 1	p_sign_posn = 0	(\$1.25)	(\$ 1.25)	(\$1.25)
	p_sign_posn = 1	+ \$1.25	+\$ 1.25	+\$1.25
	p_sign_posn = 2	\$1.25 +	\$ 1.25+	\$1.25+
	p_sign_posn = 3	+ \$1.25	+\$ 1.25	+\$1.25
p_cs_precedes = 0	p_sign_posn = 4	\$ +1.25	\$+ 1.25	+\$1.25
	p_sign_posn = 0	(1.25 \$)	(1.25 \$)	(1.25\$)
	p_sign_posn = 1	+1.25 \$	+1.25 \$	+1.25\$
	p_sign_posn = 2	1.25\$ +	1.25 \$+	1.25\$+
	p_sign_posn = 3	1.25+ \$	1.25 +\$	1.25+\$
	p_sign_posn = 4	1.25\$ +	1.25 \$+	1.25\$+

1976
1977
1978
1979

The following is an example of the interpretation of the **mon_grouping** keyword. Assuming that the value to be formatted is 123 456 789 and the **mon_thousands_sep** is ' ', then the following table shows the result. The third column shows the equivalent string in the ISO C standard that would be used by the *localeconv()* function to accommodate this grouping.

1980
1981
1982
1983
1984
1985

mon_grouping	Formatted Value	ISO C String
3;-1	123456'789	"\3\177"
3	123'456'789	"\3"
3;2;-1	1234'56'789	"\3\2\177"
3;2	12'34'56'789	"\3\2"
-1	123456789	"\177"

1986

In these examples, the octal value of {CHAR_MAX} is 177.

1987 **A.7.3.4 LC_NUMERIC**

1988

See the rationale for *LC_MONETARY* for a description of the behavior of grouping.

1989 **A.7.3.5 LC_TIME**

1990
1991
1992
1993

Although certain of the conversion specifications in the POSIX locale (such as the name of the month) are shown with initial capital letters, this need not be the case in other locales. Programs using these conversion specifications may need to adjust the capitalization if the output is going to be used at the beginning of a sentence.

1994
1995
1996

The *LC_TIME* descriptions of **abday**, **day**, **mon**, and **abmon** imply a Gregorian style calendar (7-day weeks, 12-month years, leap years, and so on). Formatting time strings for other types of calendars is outside the scope of IEEE Std 1003.1-2001.

1997
1998
1999
2000
2001
2002

While the ISO 8601:2000 standard numbers the weekdays starting with Monday, historical practice is to use the Sunday as the first day. Rather than change the order and introduce potential confusion, the days must be specified beginning with Sunday; previous references to "first day" have been removed. Note also that the Shell and Utilities volume of IEEE Std 1003.1-2001 *date* utility supports numbering compliant with the ISO 8601:2000 standard.

2003
2004
2005
2006
2007
2008

As specified under *date* in the Shell and Utilities volume of IEEE Std 1003.1-2001 and *strftime()* in the System Interfaces volume of IEEE Std 1003.1-2001, the conversion specifications corresponding to the optional keywords consist of a modifier followed by a traditional conversion specification (for instance, %Ex). If the optional keywords are not supported by the implementation or are unspecified for the current locale, these modified conversion specifications are treated as the traditional conversion specifications. For example, assume the

2009 following keywords:

2010 alt_digits "0th";"1st";"2nd";"3rd";"4th";"5th";\
2011 "6th";"7th";"8th";"9th";"10th"

2012 d_fmt "The %Od day of %B in %Y"

2013 On July 4th 1776, the %x conversion specifications would result in "The 4th day of July
2014 in 1776", while on July 14th 1789 it would result in "The 14 day of July in 1789". It
2015 can be noted that the above example is for illustrative purposes only; the %O modifier is
2016 primarily intended to provide for Kanji or Hindi digits in *date* formats.

2017 The following is an example for Japan that supports the current plus last three Emperors and
2018 reverts to Western style numbering for years prior to the Meiji era. The example also allows for
2019 the custom of using a special name for the first year of an era instead of using 1. (The examples
2020 substitute romaji where kanji should be used.)

2021 era_d_fmt "%EY%mgatsu%dnichi (%a)"

2022 era "+:2:1990/01/01:+*:Heisei:%EC%Eynen";\
2023 "+:1:1989/01/08:1989/12/31:Heisei:%ECgannen";\
2024 "+:2:1927/01/01:1989/01/07:Shouwa:%EC%Eynen";\
2025 "+:1:1926/12/25:1926/12/31:Shouwa:%ECgannen";\
2026 "+:2:1913/01/01:1926/12/24:Taishou:%EC%Eynen";\
2027 "+:1:1912/07/30:1912/12/31:Taishou:%ECgannen";\
2028 "+:2:1869/01/01:1912/07/29:Meiji:%EC%Eynen";\
2029 "+:1:1868/09/08:1868/12/31:Meiji:%ECgannen";\
2030 "-:1868:1868/09/07:-*:%Ey"

2031 Assuming that the current date is September 21, 1991, a request to *date* or *strftime()* would yield
2032 the following results:

2033 %Ec - Heisei3nen9gatsu21nichi (Sat) 14:39:26
2034 %EC - Heisei
2035 %Ex - Heisei3nen9gatsu21nichi (Sat)
2036 %Ey - 3
2037 %EY - Heisei3nen

2038 Example era definitions for the Republic of China:

2039 era "+:2:1913/01/01:+*:ChungHwaMingGuo:%EC%EyNen";\
2040 "+:1:1912/1/1:1912/12/31:ChungHwaMingGuo:%ECYuenNen";\
2041 "+:1:1911/12/31:-*:MingChien:%EC%EyNen"

2042 Example definitions for the Christian Era:

2043 era "+:1:0001/01/01:+*:AD:%EC %Ey";\
2044 "+:1:-0001/12/31:-*:BC:%Ey %EC"

2045 A.7.3.6 LC_MESSAGES

2046 The **yesstr** and **nostr** locale keywords and the YESSTR and NOSTR *langinfo* items were formerly
2047 used to match user affirmative and negative responses. In IEEE Std 1003.1-2001, the **yesexpr**,
2048 **noexpr**, YESEXPR, and NOEXPR extended regular expressions have replaced them.
2049 Applications should use the general locale-based messaging facilities to issue prompting
2050 messages which include sample desired responses.

2051 **A.7.4 Locale Definition Grammar**

2052 There is no additional rationale provided for this section.

2053 *A.7.4.1 Locale Lexical Conventions*

2054 There is no additional rationale provided for this section.

2055 *A.7.4.2 Locale Grammar*

2056 There is no additional rationale provided for this section.

2057 **A.7.5 Locale Definition Example**

2058 The following is an example of a locale definition file that could be used as input to the *localedef*
 2059 utility. It assumes that the utility is executed with the *-f* option, naming a charmap file with (at
 2060 least) the following content:

```

2061     CHARMAP
2062     <space>      \x20
2063     <dollar>     \x24
2064     <A>         \101
2065     <a>         \141
2066     <A-acute>   \346
2067     <a-acute>   \365
2068     <A-grave>   \300
2069     <a-grave>   \366
2070     <b>         \142
2071     <C>         \103
2072     <c>         \143
2073     <c-cedilla> \347
2074     <d>         \x64
2075     <H>         \110
2076     <h>         \150
2077     <eszet>    \xb7
2078     <s>         \x73
2079     <z>         \x7a
2080     END CHARMAP

```

2081 It should not be taken as complete or to represent any actual locale, but only to illustrate the
 2082 syntax.

```

2083     #
2084     LC_CTYPE
2085     lower  <a>;<b>;<c>;<c-cedilla>;<d>;...;<z>
2086     upper  A;B;C;Ç;...;Z
2087     space  \x20;\x09;\x0a;\x0b;\x0c;\x0d
2088     blank  \040;\011
2089     toupper (<a>,<A>);(b,B);(c,C);(ç,Ç);(d,D);(z,Z)
2090     END LC_CTYPE
2091     #
2092     LC_COLLATE
2093     #
2094     # The following example of collation is based on
2095     # Canadian standard Z243.4.1-1998, "Canadian Alphanumeric
2096     # Ordering Standard for Character Sets of CSA Z234.4 Standard".

```

```
2097 # (Other parts of this example locale definition file do not
2098 # purport to relate to Canada, or to any other real culture.)
2099 # The proposed standard defines a 4-weight collation, such that
2100 # in the first pass, characters are compared without regard to
2101 # case or accents; in the second pass, backwards-compare without
2102 # regard to case; in the third pass, forwards-compare without
2103 # regard to diacriticals. In the 3 first passes, non-alphabetic
2104 # characters are ignored; in the fourth pass, only special
2105 # characters are considered, such that "The string that has a
2106 # special character in the lowest position comes first. If two
2107 # strings have a special character in the same position, the
2108 # collation value of the special character determines ordering.
2109 #
2110 # Only a subset of the character set is used here; mostly to
2111 # illustrate the set-up.
2112 #
2113 collating-symbol <NULL>
2114 collating-symbol <LOW_VALUE>
2115 collating-symbol <LOWER-CASE>
2116 collating-symbol <SUBSCRIPT-LOWER>
2117 collating-symbol <SUPERSCRIPT-LOWER>
2118 collating-symbol <UPPER-CASE>
2119 collating-symbol <NO-ACCENT>
2120 collating-symbol <PECULIAR>
2121 collating-symbol <LIGATURE>
2122 collating-symbol <ACUTE>
2123 collating-symbol <GRAVE>
2124 # Further collating-symbols follow.
2125 #
2126 # Properly, the standard does not include any multi-character
2127 # collating elements; the one below is added for completeness.
2128 #
2129 collating_element <ch> from "<c><h>"
2130 collating_element <CH> from "<C><H>"
2131 collating_element <Ch> from "<C><h>"
2132 #
2133 order_start forward;backward;forward;forward,position
2134 #
2135 # Collating symbols are specified first in the sequence to allocate
2136 # basic collation values to them, lower than that of any character.
2137 <NULL>
2138 <LOW_VALUE>
2139 <LOWER-CASE>
2140 <SUBSCRIPT-LOWER>
2141 <SUPERSCRIPT-LOWER>
2142 <UPPER-CASE>
2143 <NO-ACCENT>
2144 <PECULIAR>
2145 <LIGATURE>
2146 <ACUTE>
2147 <GRAVE>
2148 <RING-ABOVE>
```

```

2149 <DIAERESIS>
2150 <TILDE>
2151 # Further collating symbols are given a basic collating value here.
2152 #
2153 # Here follow special characters.
2154 <space>          IGNORE;IGNORE;IGNORE;<space>
2155 # Other special characters follow here.
2156 #
2157 # Here follow the regular characters.
2158 <a>              <a>;<NO-ACCENT>;<LOWER-CASE>;IGNORE
2159 <A>              <a>;<NO-ACCENT>;<UPPER-CASE>;IGNORE
2160 <a-acute>        <a>;<ACUTE>;<LOWER-CASE>;IGNORE
2161 <A-acute>        <a>;<ACUTE>;<UPPER-CASE>;IGNORE
2162 <a-grave>        <a>;<GRAVE>;<LOWER-CASE>;IGNORE
2163 <A-grave>        <a>;<GRAVE>;<UPPER-CASE>;IGNORE
2164 <ae>            "<a><e>";"<LIGATURE><LIGATURE>";\
2165                "<LOWER-CASE><LOWER-CASE>";IGNORE
2166 <AE>            "<a><e>";"<LIGATURE><LIGATURE>";\
2167                "<UPPER-CASE><UPPER-CASE>";IGNORE
2168 <b>              <b>;<NO-ACCENT>;<LOWER-CASE>;IGNORE
2169 <B>              <b>;<NO-ACCENT>;<UPPER-CASE>;IGNORE
2170 <c>              <c>;<NO-ACCENT>;<LOWER-CASE>;IGNORE
2171 <C>              <c>;<NO-ACCENT>;<UPPER-CASE>;IGNORE
2172 <ch>            <ch>;<NO-ACCENT>;<LOWER-CASE>;IGNORE
2173 <Ch>            <ch>;<NO-ACCENT>;<PECULIAR>;IGNORE
2174 <CH>            <ch>;<NO-ACCENT>;<UPPER-CASE>;IGNORE
2175 #
2176 # As an example, the strings "Bach" and "bach" could be encoded (for
2177 # compare purposes) as:
2178 # "Bach" <b>;<a>;<ch>;<LOW_VALUE>;<NO-ACCENT>;<NO-ACCENT>;\
2179 #        <NO-ACCENT>;<LOW_VALUE>;<UPPER-CASE>;<LOWER-CASE>;\
2180 #        <LOWER-CASE>;<NULL>
2181 # "bach" <b>;<a>;<ch>;<LOW_VALUE>;<NO-ACCENT>;<NO-ACCENT>;\
2182 #        <NO-ACCENT>;<LOW_VALUE>;<LOWER-CASE>;<LOWER-CASE>;\
2183 #        <LOWER-CASE>;<NULL>
2184 #
2185 # The two strings are equal in pass 1 and 2, but differ in pass 3.
2186 #
2187 # Further characters follow.
2188 #
2189 UNDEFINED      IGNORE;IGNORE;IGNORE;IGNORE
2190 #
2191 order_end
2192 #
2193 END LC_COLLATE
2194 #
2195 LC_MONETARY
2196 int_curr_symbol    "USD "
2197 currency_symbol    "$"
2198 mon_decimal_point  "."
2199 mon_grouping       3;0
2200 positive_sign      ""

```



```

2201     negative_sign      "-"
2202     p_cs_precedes      1
2203     n_sign_posn        0
2204     END LC_MONETARY
2205     #
2206     LC_NUMERIC
2207     copy "US_en.ASCII"
2208     END LC_NUMERIC
2209     #
2210     LC_TIME
2211     abday  "Sun"; "Mon"; "Tue"; "Wed"; "Thu"; "Fri"; "Sat"
2212     #
2213     day    "Sunday"; "Monday"; "Tuesday"; "Wednesday"; \
2214           "Thursday"; "Friday"; "Saturday"
2215     #
2216     abmon  "Jan"; "Feb"; "Mar"; "Apr"; "May"; "Jun"; \
2217           "Jul"; "Aug"; "Sep"; "Oct"; "Nov"; "Dec"
2218     #
2219     mon    "January"; "February"; "March"; "April"; \
2220           "May"; "June"; "July"; "August"; "September"; \
2221           "October"; "November"; "December"
2222     #
2223     d_t_fmt "%a %b %d %T %Z %Y\n"
2224     END LC_TIME
2225     #
2226     LC_MESSAGES
2227     yesexpr  "^[yY][[:alpha:]]*" | "(OK)"
2228     #
2229     noexpr   "^[nN][[:alpha:]]*"
2230     END LC_MESSAGES

```

2231 A.8 Environment Variables

2232 A.8.1 Environment Variable Definition

2233 The variable *environ* is not intended to be declared in any header, but rather to be declared by the
 2234 user for accessing the array of strings that is the environment. This is the traditional usage of the
 2235 symbol. Putting it into a header could break some programs that use the symbol for their own
 2236 purposes.

2237 The decision to restrict conforming systems to the use of digits, uppercase letters, and
 2238 underscores for environment variable names allows applications to use lowercase letters in their
 2239 environment variable names without conflicting with any conforming system.

2240 In addition to the obvious conflict with the shell syntax for positional parameter substitution,
 2241 some historical applications (including some shells) exclude names with leading digits from the
 2242 environment.

2243 **A.8.2 Internationalization Variables**

2244 The text about locale implies that any utilities written in standard C and conforming to
2245 IEEE Std 1003.1-2001 must issue the following call:

```
2246     setlocale(LC_ALL, " ")
```

2247 If this were omitted, the ISO C standard specifies that the C locale would be used.

2248 If any of the environment variables are invalid, it makes sense to default to an implementation-
2249 defined, consistent locale environment. It is more confusing for a user to have partial settings
2250 occur in case of a mistake. All utilities would then behave in one language/cultural
2251 environment. Furthermore, it provides a way of forcing the whole environment to be the
2252 implementation-defined default. Disastrous results could occur if a pipeline of utilities partially
2253 uses the environment variables in different ways. In this case, it would be appropriate for
2254 utilities that use *LANG* and related variables to exit with an error if any of the variables are
2255 invalid. For example, users typing individual commands at a terminal might want *date* to work if
2256 *LC_MONETARY* is invalid as long as *LC_TIME* is valid. Since these are conflicting reasonable
2257 alternatives, IEEE Std 1003.1-2001 leaves the results unspecified if the locale environment
2258 variables would not produce a complete locale matching the specification of the user.

2259 The locale settings of individual categories cannot be truly independent and still guarantee
2260 correct results. For example, when collating two strings, characters must first be extracted from
2261 each string (governed by *LC_CTYPE*) before being mapped to collating elements (governed by
2262 *LC_COLLATE*) for comparison. That is, if *LC_CTYPE* is causing parsing according to the rules of
2263 a large, multi-byte code set (potentially returning 20 000 or more distinct character codeset
2264 values), but *LC_COLLATE* is set to handle only an 8-bit codeset with 256 distinct characters,
2265 meaningful results are obviously impossible.

2266 The *LC_MESSAGES* variable affects the language of messages generated by the standard
2267 utilities.

2268 The description of the environment variable names starting with the characters “LC_”
2269 acknowledges the fact that the interfaces presented may be extended as new international
2270 functionality is required. In the ISO C standard, names preceded by “LC_” are reserved in the
2271 name space for future categories.

2272 To avoid name clashes, new categories and environment variables are divided into two
2273 classifications: “implementation-independent” and “implementation-defined”.

2274 Implementation-independent names will have the following format:

```
2275     LC_NAME
```

2276 where *NAME* is the name of the new category and environment variable. Capital letters must be
2277 used for implementation-independent names.

2278 Implementation-defined names must be in lowercase letters, as below:

```
2279     lc_name
```

2280 **A.8.3 Other Environment Variables**2281 **COLUMNS, LINES**

2282 The default values for the number of column positions, *COLUMNS*, and screen height, *LINES*,
2283 are unspecified because historical implementations use different methods to determine values
2284 corresponding to the size of the screen in which the utility is run. This size is typically known to
2285 the implementation through the value of *TERM*, or by more elaborate methods such as
2286 extensions to the *stty* utility or knowledge of how the user is dynamically resizing windows on a
2287 bit-mapped display terminal. Users should not need to set these variables in the environment
2288 unless there is a specific reason to override the default behavior of the implementation, such as
2289 to display data in an area arbitrarily smaller than the terminal or window. Values for these
2290 variables that are not decimal integers greater than zero are implicitly undefined values; it is
2291 unnecessary to enumerate all of the possible values outside of the acceptable set.

2292 **LOGNAME**

2293 In most implementations, the value of such a variable is easily forged, so security-critical
2294 applications should rely on other means of determining user identity. *LOGNAME* is required to
2295 be constructed from the portable filename character set for reasons of interchange. No diagnostic
2296 condition is specified for violating this rule, and no requirement for enforcement exists. The
2297 intent of the requirement is that if extended characters are used, the “guarantee” of portability
2298 implied by a standard is void.

2299 **PATH**

2300 Many historical implementations of the Bourne shell do not interpret a trailing colon to represent
2301 the current working directory and are thus non-conforming. The C Shell and the KornShell
2302 conform to IEEE Std 1003.1-2001 on this point. The usual name of dot may also be used to refer
2303 to the current working directory.

2304 Many implementations historically have used a default value of */bin* and */usr/bin* for the *PATH*
2305 variable. IEEE Std 1003.1-2001 does not mandate this default path be identical to that retrieved
2306 from *getconf _CS_PATH* because it is likely that the standardized utilities may be provided in
2307 another directory separate from the directories used by some historical applications.

2308 **SHELL**

2309 The *SHELL* variable names the preferred shell of the user; it is a guide to applications. There is
2310 no direct requirement that that shell conform to IEEE Std 1003.1-2001; that decision should rest
2311 with the user. It is the intention of the standard developers that alternative shells be permitted, if
2312 the user chooses to develop or acquire one. An operating system that builds its shell into the
2313 “kernel” in such a manner that alternative shells would be impossible does not conform to the
2314 spirit of IEEE Std 1003.1-2001.

2315 **TZ**

2316 The quoted form of the timezone variable allows timezone names of the form UTC+1 (or any
2317 name that contains the character plus ('+'), the character minus ('-'), or digits), which may be
2318 appropriate for countries that do not have an official timezone name. It would be coded as
2319 <UTC+1>+1<UTC+2>, which would cause *std* to have a value of UTC+1 and *dst* a value of
2320 UTC+2, each with a length of 5 characters. This does not appear to conflict with any existing
2321 usage. The characters '<' and '>' were chosen for quoting because they are easier to parse
2322 visually than a quoting character that does not provide some sense of bracketing (and in a string
2323 like this, such bracketing is helpful). They were also chosen because they do not need special

2324 treatment when assigning to the *TZ* variable. Users are often confused by embedding quotes in a
 2325 string. Because '*<*' and '*>*' are meaningful to the shell, the whole string would have to be
 2326 quoted, but that is easily explained. (Parentheses would have presented the same problems.)
 2327 Although the '*>*' symbol could have been permitted in the string by either escaping it or
 2328 doubling it, it seemed of little value to require that. This could be provided as an extension if
 2329 there was a need. Timezone names of this new form lead to a requirement that the value of
 2330 `{_POSIX_TZNAME_MAX}` change from 3 to 6.

2331 Since the *TZ* environment variable is usually inherited by all applications started by a user after
 2332 the value of the *TZ* environment variable is changed and since many applications run using the
 2333 C or POSIX locale, using characters that are not in the portable character set in the *std* and *dst*
 2334 fields could cause unexpected results.

2335 The format of the *TZ* environment variable is changed in Issue 6 to allow for the quoted form, as
 2336 defined in previous versions of the ISO POSIX-1 standard.

2337 **A.9 Regular Expressions**

2338 Rather than repeating the description of REs for each utility supporting REs, the standard
 2339 developers preferred a common, comprehensive description of regular expressions in one place.
 2340 The most common behavior is described here, and exceptions or extensions to this are
 2341 documented for the respective utilities, as appropriate.

2342 The BRE corresponds to the *ed* or historical *grep* type, and the ERE corresponds to the historical
 2343 *egrep* type (now *grep -E*).

2344 The text is based on the *ed* description and substantially modified, primarily to aid developers
 2345 and others in the understanding of the capabilities and limitations of REs. Much of this was
 2346 influenced by internationalization requirements.

2347 It should be noted that the definitions in this section do not cover the *tr* utility; the *tr* syntax does
 2348 not employ REs.

2349 The specification of REs is particularly important to internationalization because pattern
 2350 matching operations are very basic operations in business and other operations. The syntax and
 2351 rules of REs are intended to be as intuitive as possible to make them easy to understand and use.
 2352 The historical rules and behavior do not provide that capability to non-English language users,
 2353 and do not provide the necessary support for commonly used characters and language
 2354 constructs. It was necessary to provide extensions to the historical RE syntax and rules to
 2355 accommodate other languages.

2356 As they are limited to bracket expressions, the rationale for these modifications is in the Base
 2357 Definitions volume of IEEE Std 1003.1-2001, Section 9.3.5, RE Bracket Expression.

2358 **A.9.1 Regular Expression Definitions**

2359 It is possible to determine what strings correspond to subexpressions by recursively applying
 2360 the leftmost longest rule to each subexpression, but only with the proviso that the overall match
 2361 is leftmost longest. For example, matching "`\(ac*\\)c*d[ac]*\1`" against *acdacaaa* matches
 2362 *acdacaaa* (with `\1=a`); simply matching the longest match for "`\(ac*\\)`" would yield `\1=ac`, but
 2363 the overall match would be smaller (*acdac*). Conceptually, the implementation must examine
 2364 every possible match and among those that yield the leftmost longest total matches, pick the one
 2365 that does the longest match for the leftmost subexpression, and so on. Note that this means that
 2366 matching by subexpressions is context-dependent: a subexpression within a larger RE may
 2367 match a different string from the one it would match as an independent RE, and two instances of

2368 the same subexpression within the same larger RE may match different lengths even in similar
 2369 sequences of characters. For example, in the ERE "(a.*b)(a.*b)", the two identical
 2370 subexpressions would match four and six characters, respectively, of *accbacccb*.

2371 The definition of single character has been expanded to include also collating elements
 2372 consisting of two or more characters; this expansion is applicable only when a bracket
 2373 expression is included in the BRE or ERE. An example of such a collating element may be the
 2374 Dutch *ij*, which collates as a 'y'. In some encodings, a ligature "i with j" exists as a character
 2375 and would represent a single-character collating element. In another encoding, no such ligature
 2376 exists, and the two-character sequence *ij* is defined as a multi-character collating element.
 2377 Outside brackets, the *ij* is treated as a two-character RE and matches the same characters in a
 2378 string. Historically, a bracket expression only matched a single character. The ISO POSIX-2:1993
 2379 standard required bracket expressions like "[^[:lower:]]" to match multi-character collating
 2380 elements such as "ij". However, this requirement led to behavior that many users did not
 2381 expect and that could not feasibly be mimicked in user code, and it was rarely if ever
 2382 implemented correctly. The current standard leaves it unspecified whether a bracket expression
 2383 matches a multi-character collating element, allowing both historical and ISO POSIX-2:1993
 2384 standard implementations to conform.

2385 Also, in the current standard, it is unspecified whether character class expressions like
 2386 "[[:lower:]]" can include multi-character collating elements like "ij"; hence
 2387 "[[:lower:]]" can match "ij", and "[^[:lower:]]" can fail to match "ij". Common
 2388 practice is for a character class expression to match a collating element if it matches the collating
 2389 element's first character.

2390 A.9.2 Regular Expression General Requirements

2391 The definition of which sequence is matched when several are possible is based on the leftmost-
 2392 longest rule historically used by deterministic recognizers. This rule is easier to define and
 2393 describe, and arguably more useful, than the first-match rule historically used by non-
 2394 deterministic recognizers. It is thought that dependencies on the choice of rule are rare; carefully
 2395 contrived examples are needed to demonstrate the difference.

2396 A formal expression of the leftmost-longest rule is:

2397 The search is performed as if all possible suffixes of the string were tested for a prefix
 2398 matching the pattern; the longest suffix containing a matching prefix is chosen, and the
 2399 longest possible matching prefix of the chosen suffix is identified as the matching sequence.

2400 Historically, most RE implementations only match lines, not strings. However, that is more an
 2401 effect of the usage than of an inherent feature of REs themselves. Consequently,
 2402 IEEE Std 1003.1-2001 does not regard <newline>s as special; they are ordinary characters, and
 2403 both a period and a non-matching list can match them. Those utilities (like *grep*) that do not
 2404 allow <newline>s to match are responsible for eliminating any <newline> from strings before
 2405 matching against the RE. The *regcomp()* function, however, can provide support for such
 2406 processing without violating the rules of this section.

2407 Some implementations of *egrep* have had very limited flexibility in handling complex EREs.
 2408 IEEE Std 1003.1-2001 does not attempt to define the complexity of a BRE or ERE, but does place a
 2409 lower limit on it—any RE must be handled, as long as it can be expressed in 256 bytes or less. (Of
 2410 course, this does not place an upper limit on the implementation.) There are historical programs
 2411 using a non-deterministic-recognizer implementation that should have no difficulty with this
 2412 limit. It is possible that a good approach would be to attempt to use the faster, but more limited,
 2413 deterministic recognizer for simple expressions and to fall back on the non-deterministic
 2414 recognizer for those expressions requiring it. Non-deterministic implementations must be
 2415 careful to observe the rules on which match is chosen; the longest match, not the first match,

2416 starting at a given character is used.

2417 The term “invalid” highlights a difference between this section and some others:
 2418 IEEE Std 1003.1-2001 frequently avoids mandating of errors for syntax violations because they
 2419 can be used by implementors to trigger extensions. However, the authors of the
 2420 internationalization features of REs wanted to mandate errors for certain conditions to identify
 2421 usage problems or non-portable constructs. These are identified within this rationale as
 2422 appropriate. The remaining syntax violations have been left implicitly or explicitly undefined.
 2423 For example, the BRE construct " $\{1,2,3\}$ " does not comply with the grammar. A
 2424 conforming application cannot rely on it producing an error nor matching the literal characters
 2425 " $\{1,2,3\}$ ". The term “undefined” was used in favor of “unspecified” because many of the
 2426 situations are considered errors on some implementations, and the standard developers
 2427 considered that consistency throughout the section was preferable to mixing undefined and
 2428 unspecified.

2429 **A.9.3 Basic Regular Expressions**

2430 There is no additional rationale provided for this section.

2431 *A.9.3.1 BREs Matching a Single Character or Collating Element*

2432 There is no additional rationale provided for this section.

2433 *A.9.3.2 BRE Ordinary Characters*

2434 There is no additional rationale provided for this section.

2435 *A.9.3.3 BRE Special Characters*

2436 There is no additional rationale provided for this section.

2437 *A.9.3.4 Periods in BREs*

2438 There is no additional rationale provided for this section.

2439 *A.9.3.5 RE Bracket Expression*

2440 Range expressions are, historically, an integral part of REs. However, the requirements of
 2441 “natural language behavior” and portability do conflict. In the POSIX locale, ranges must be
 2442 treated according to the collating sequence and include such characters that fall within the range
 2443 based on that collating sequence, regardless of character values. In other locales, ranges have
 2444 unspecified behavior.

2445 Some historical implementations allow range expressions where the ending range point of one
 2446 range is also the starting point of the next (for instance, " $[a-m-o]$ "). This behavior should not
 2447 be permitted, but to avoid breaking historical implementations, it is now *undefined* whether it is a
 2448 valid expression and how it should be interpreted.

2449 Current practice in *awk* and *lex* is to accept escape sequences in bracket expressions as per the
 2450 Base Definitions volume of IEEE Std 1003.1-2001, Table 5-1, Escape Sequences and Associated
 2451 Actions, while the normal ERE behavior is to regard such a sequence as consisting of two
 2452 characters. Allowing the *awk/lex* behavior in EREs would change the normal behavior in an
 2453 unacceptable way; it is expected that *awk* and *lex* will decode escape sequences in EREs before
 2454 passing them to *regcomp()* or comparable routines. Each utility describes the escape sequences it
 2455 accepts as an exception to the rules in this section; the list is not the same, for historical reasons.

2456 As noted previously, the new syntax and rules have been added to accommodate other
 2457 languages than English. The remainder of this section describes the rationale for these
 2458 modifications.

2459 In the POSIX locale, a regular expression that starts with a range expression matches a set of
 2460 strings that are contiguously sorted, but this is not necessarily true in other locales. For example,
 2461 a French locale might have the following behavior:

```
2462 $ ls
2463 alpha Alpha estimé ESTIMÉ été eurêka
2464 $ ls [a-e]*
2465 alpha Alpha estimé eurêka
```

2466 Such disagreements between matching and contiguous sorting are unavoidable because POSIX
 2467 sorting cannot be implemented in terms of a deterministic finite-state automaton (DFA), but
 2468 range expressions by design are implementable in terms of DFAs.

2469 Historical implementations used native character order to interpret range expressions. The
 2470 ISO POSIX-2:1993 standard instead required collating element order (CEO): the order that
 2471 collating elements were specified between the **order_start** and **order_end** keywords in the
 2472 *LC_COLLATE* category of the current locale. CEO had some advantages in portability over the
 2473 native character order, but it also had some disadvantages:

- 2474 • CEO could not feasibly be mimicked in user code, leading to inconsistencies between POSIX
 2475 matchers and matchers in popular user programs like Emacs, *ksh*, and Perl.
- 2476 • CEO caused range expressions to match accented and capitalized letters contrary to many
 2477 users' expectations. For example, "[a-e]" typically matched both 'E' and 'â' but neither
 2478 'A' nor 'é'.
- 2479 • CEO was not consistent across implementations. In practice, CEO was often less portable
 2480 than native character order. For example, it was common for the CEOs of two
 2481 implementation-supplied locales to disagree, even if both locales were named "da_DK".

2482 Because of these problems, some implementations of regular expressions continued to use
 2483 native character order. Others used the collation sequence, which is more consistent with sorting
 2484 than either CEO or native order, but which departs further from the traditional POSIX semantics
 2485 because it generally requires "[a-e]" to match either 'A' or 'E' but not both. As a result of
 2486 this kind of implementation variation, programmers who wanted to write portable regular
 2487 expressions could not rely on the ISO POSIX-2:1993 standard guarantees in practice.

2488 While revising the standard, lengthy consideration was given to proposals to attack this problem
 2489 by adding an API for querying the CEO to allow user-mode matchers, but none of these
 2490 proposals had implementation experience and none achieved consensus. Leaving the standard
 2491 alone was also considered, but rejected due to the problems described above.

2492 The current standard leaves unspecified the behavior of a range expression outside the POSIX
 2493 locale. This makes it clearer that conforming applications should avoid range expressions
 2494 outside the POSIX locale, and it allows implementations and compatible user-mode matchers to
 2495 interpret range expressions using native order, CEO, collation sequence, or other, more
 2496 advanced techniques. The concerns which led to this change were raised in IEEE PASC
 2497 interpretation 1003.2 #43 and others, and related to ambiguities in the specification of how
 2498 multi-character collating elements should be handled in range expressions. These ambiguities
 2499 had led to multiple interpretations of the specification, in conflicting ways, which led to varying
 2500 implementations. As noted above, efforts were made to resolve the differences, but no solution
 2501 has been found that would be specific enough to allow for portable software while not
 2502 invalidating existing implementations.

2503 The standard developers recognize that collating elements are important, such elements being
 2504 common in several European languages; for example, 'ch' or 'll' in traditional Spanish; 'aa'
 2505 in several Scandinavian languages. Existing internationalized implementations have processed,
 2506 and continue to process, these elements in range expressions. Efforts are expected to continue in
 2507 the future to find a way to define the behavior of these elements precisely and portably.

2508 The ISO POSIX-2:1993 standard required "[b-a]" to be an invalid expression in the POSIX
 2509 locale, but this requirement has been relaxed in this version of the standard so that "[b-a]" can
 2510 instead be treated as a valid expression that does not match any string.

2511 A.9.3.6 BREs Matching Multiple Characters

2512 The limit of nine back-references to subexpressions in the RE is based on the use of a single-digit
 2513 identifier; increasing this to multiple digits would break historical applications. This does not
 2514 imply that only nine subexpressions are allowed in REs. The following is a valid BRE with ten
 2515 subexpressions:

```
2516 \(\(\(ab\) *c\) *d\)\(ef\) *(gh)\{2\}\(ij\) *(kl\) *(mn\) *(op\) *(qr\) *
```

2517 The standard developers regarded the common historical behavior, which supported "\n*", but
 2518 not "\n{min,max}", "\(...)*", or "\(...)\{min,max\}", as a non-intentional
 2519 result of a specific implementation, and they supported both duplication and interval
 2520 expressions following subexpressions and back-references.

2521 The changes to the processing of the back-reference expression remove an unspecified or
 2522 ambiguous behavior in the Shell and Utilities volume of IEEE Std 1003.1-2001, aligning it with
 2523 the requirements specified for the *regcomp()* expression, and is the result of PASC Interpretation
 2524 1003.2-92 #43 submitted for the ISO POSIX-2:1993 standard.

2525 A.9.3.7 BRE Precedence

2526 There is no additional rationale provided for this section.

2527 A.9.3.8 BRE Expression Anchoring

2528 Often, the dollar sign is viewed as matching the ending <newline> in text files. This is not
 2529 strictly true; the <newline> is typically eliminated from the strings to be matched, and the dollar
 2530 sign matches the terminating null character.

2531 The ability of '^', '\$', and '*' to be non-special in certain circumstances may be confusing to
 2532 some programmers, but this situation was changed only in a minor way from historical practice
 2533 to avoid breaking many historical scripts. Some consideration was given to making the use of
 2534 the anchoring characters undefined if not escaped and not at the beginning or end of strings.
 2535 This would cause a number of historical BREs, such as "2^10", "\$HOME", and "\$1.35", that
 2536 relied on the characters being treated literally, to become invalid.

2537 However, one relatively uncommon case was changed to allow an extension used on some
 2538 implementations. Historically, the BREs "^foo" and "\(^foo)" did not match the same
 2539 string, despite the general rule that subexpressions and entire BREs match the same strings. To
 2540 increase consensus, IEEE Std 1003.1-2001 has allowed an extension on some implementations to
 2541 treat these two cases in the same way by declaring that anchoring *may* occur at the beginning or
 2542 end of a subexpression. Therefore, portable BREs that require a literal circumflex at the
 2543 beginning or a dollar sign at the end of a subexpression must escape them. Note that a BRE such
 2544 as "a\(^bc)" will either match "a^bc" or nothing on different systems under the rules.

2545 ERE anchoring has been different from BRE anchoring in all historical systems. An unescaped
 2546 anchor character has never matched its literal counterpart outside a bracket expression. Some

2547 implementations treated "foo\$bar" as a valid expression that never matched anything; others
 2548 treated it as invalid. IEEE Std 1003.1-2001 mandates the former, valid unmatched behavior.

2549 Some implementations have extended the BRE syntax to add alternation. For example, the
 2550 subexpression "\ (foo\$|bar\)" would match either "foo" at the end of the string or "bar"
 2551 anywhere. The extension is triggered by the use of the undefined "\|" sequence. Because the
 2552 BRE is undefined for portable scripts, the extending system is free to make other assumptions,
 2553 such that the '\$' represents the end-of-line anchor in the middle of a subexpression. If it were
 2554 not for the extension, the '\$' would match a literal dollar sign under the rules.

2555 **A.9.4 Extended Regular Expressions**

2556 As with BREs, the standard developers decided to make the interpretation of escaped ordinary
 2557 characters undefined.

2558 The right parenthesis is not listed as an ERE special character because it is only special in the
 2559 context of a preceding left parenthesis. If found without a preceding left parenthesis, the right
 2560 parenthesis has no special meaning.

2561 The interval expression, "{m,n}", has been added to EREs. Historically, the interval expression
 2562 has only been supported in some ERE implementations. The standard developers estimated that
 2563 the addition of interval expressions to EREs would not decrease consensus and would also make
 2564 BREs more of a subset of EREs than in many historical implementations.

2565 It was suggested that, in addition to interval expressions, back-references ('\n') should also be
 2566 added to EREs. This was rejected by the standard developers as likely to decrease consensus.

2567 In historical implementations, multiple duplication symbols are usually interpreted from left to
 2568 right and treated as additive. As an example, "a+b" matches zero or more instances of 'a'
 2569 followed by a 'b'. In IEEE Std 1003.1-2001, multiple duplication symbols are undefined; that is,
 2570 they cannot be relied upon for conforming applications. One reason for this is to provide some
 2571 scope for future enhancements.

2572 The precedence of operations differs between EREs and those in *lex*; in *lex*, for historical reasons,
 2573 interval expressions have a lower precedence than concatenation.

2574 *A.9.4.1 EREs Matching a Single Character or Collating Element*

2575 There is no additional rationale provided for this section.

2576 *A.9.4.2 ERE Ordinary Characters*

2577 There is no additional rationale provided for this section.

2578 *A.9.4.3 ERE Special Characters*

2579 There is no additional rationale provided for this section.

2580 *A.9.4.4 Periods in EREs*

2581 There is no additional rationale provided for this section.

2582 A.9.4.5 *ERE Bracket Expression*

2583 There is no additional rationale provided for this section.

2584 A.9.4.6 *EREs Matching Multiple Characters*

2585 There is no additional rationale provided for this section.

2586 A.9.4.7 *ERE Alternation*

2587 There is no additional rationale provided for this section.

2588 A.9.4.8 *ERE Precedence*

2589 There is no additional rationale provided for this section.

2590 A.9.4.9 *ERE Expression Anchoring*

2591 There is no additional rationale provided for this section.

2592 **A.9.5 Regular Expression Grammar**

2593 The grammars are intended to represent the range of acceptable syntaxes available to
 2594 conforming applications. There are instances in the text where undefined constructs are
 2595 described; as explained previously, these allow implementation extensions. There is no intended
 2596 requirement that an implementation extension must somehow fit into the grammars shown
 2597 here.

2598 The BRE grammar does not permit L_ANCHOR or R_ANCHOR inside "\(" and "\)" (which
 2599 implies that '^' and '\$' are ordinary characters). This reflects the semantic limits on the
 2600 application, as noted in the Base Definitions volume of IEEE Std 1003.1-2001, Section 9.3.8, BRE
 2601 Expression Anchoring. Implementations are permitted to extend the language to interpret '^'
 2602 and '\$' as anchors in these locations, and as such, conforming applications cannot use
 2603 unescaped '^' and '\$' in positions inside "\(" and "\)" that might be interpreted as anchors.

2604 The ERE grammar does not permit several constructs that the Base Definitions volume of
 2605 IEEE Std 1003.1-2001, Section 9.4.2, ERE Ordinary Characters and the Base Definitions volume of
 2606 IEEE Std 1003.1-2001, Section 9.4.3, ERE Special Characters specify as having undefined results:

- 2607 • ORD_CHAR preceded by '\'
- 2608 • *ERE_dupl_symbol*(s) appearing first in an ERE, or immediately following '|', '^', or '('
- 2609 • '{' not part of a valid *ERE_dupl_symbol*
- 2610 • '|' appearing first or last in an ERE, or immediately following '|' or '(', or immediately
 2611 preceding ')'

2612 Implementations are permitted to extend the language to allow these. Conforming applications
 2613 cannot use such constructs.

2614 A.9.5.1 *BRE/ERE Grammar Lexical Conventions*

2615 There is no additional rationale provided for this section.

2616 A.9.5.2 *RE and Bracket Expression Grammar*

2617 The removal of the *Back_open_paren Back_close_paren* option from the *nondupl_RE* specification is
2618 the result of PASC Interpretation 1003.2-92 #43 submitted for the ISO POSIX-2: 1993 standard.
2619 Although the grammar required support for null subexpressions, this section does not describe
2620 the meaning of, and historical practice did not support, this construct.

2621 A.9.5.3 *ERE Grammar*

2622 There is no additional rationale provided for this section.

2623 **A.10 Directory Structure and Devices**2624 **A.10.1 Directory Structure and Files**

2625 A description of the historical */usr/tmp* was omitted, removing any concept of differences in
2626 emphasis between the */* and */usr* directories. The descriptions of */bin*, */usr/bin*, */lib*, and */usr/lib*
2627 were omitted because they are not useful for applications. In an early draft, a distinction was
2628 made between system and application directory usage, but this was not found to be useful.

2629 The directories */* and */dev* are included because the notion of a hierarchical directory structure is
2630 key to other information presented elsewhere in IEEE Std 1003.1-2001. In early drafts, it was
2631 argued that special devices and temporary files could conceivably be handled without a
2632 directory structure on some implementations. For example, the system could treat the characters
2633 "*/tmp*" as a special token that would store files using some non-POSIX file system structure.
2634 This notion was rejected by the standard developers, who required that all the files in this
2635 section be implemented via POSIX file systems.

2636 The */tmp* directory is retained in IEEE Std 1003.1-2001 to accommodate historical applications
2637 that assume its availability. Implementations are encouraged to provide suitable directory
2638 names in the environment variable *TMPDIR* and applications are encouraged to use the contents
2639 of *TMPDIR* for creating temporary files.

2640 The standard files */dev/null* and */dev/tty* are required to be both readable and writable to allow
2641 applications to have the intended historical access to these files.

2642 The standard file */dev/console* has been added for alignment with the Single UNIX Specification.

2643 **A.10.2 Output Devices and Terminal Types**

2644 There is no additional rationale provided for this section.

2645 **A.11 General Terminal Interface**

2646 If the implementation does not support this interface on any device types, it should behave as if
2647 it were being used on a device that is not a terminal device (in most cases *errno* will be set to
2648 [ENOTTY] on return from functions defined by this interface). This is based on the fact that
2649 many applications are written to run both interactively and in some non-interactive mode, and
2650 they adapt themselves at runtime. Requiring that they all be modified to test an environment
2651 variable to determine whether they should try to adapt is unnecessary. On a system that
2652 provides no general terminal interface, providing all the entry points as stubs that return
2653 [ENOTTY] (or an equivalent, as appropriate) has the same effect and requires no changes to the
2654 application.

2655 Although the needs of both interface implementors and application developers were addressed
2656 throughout IEEE Std 1003.1-2001, this section pays more attention to the needs of the latter. This
2657 is because, while many aspects of the programming interface can be hidden from the user by the
2658 application developer, the terminal interface is usually a large part of the user interface.
2659 Although to some extent the application developer can build missing features or work around
2660 inappropriate ones, the difficulties of doing that are greater in the terminal interface than
2661 elsewhere. For example, efficiency prohibits the average program from interpreting every
2662 character passing through it in order to simulate character erase, line kill, and so on. These
2663 functions should usually be done by the operating system, possibly at the interrupt level.

2664 The *tc**() functions were introduced as a way of avoiding the problems inherent in the
2665 traditional *ioctl*() function and in variants of it that were proposed. For example, *tcsetattr*() is
2666 specified in place of the use of the TCSETA *ioctl*() command function. This allows specification
2667 of all the arguments in a manner consistent with the ISO C standard unlike the varying third
2668 argument of *ioctl*(), which is sometimes a pointer (to any of many different types) and
2669 sometimes an **int**.

2670 The advantages of this new method include:

- 2671 • It allows strict type checking.
- 2672 • The direction of transfer of control data is explicit.
- 2673 • Portable capabilities are clearly identified.
- 2674 • The need for a general interface routine is avoided.
- 2675 • Size of the argument is well-defined (there is only one type).

2676 The disadvantages include:

- 2677 • No historical implementation used the new method.
- 2678 • There are many small routines instead of one general-purpose one.
- 2679 • The historical parallel with *fcntl*() is broken.

2680 The issue of modem control was excluded from IEEE Std 1003.1-2001 on the grounds that:

- 2681 • It was concerned with setting and control of hardware timers.
- 2682 • The appropriate timers and settings vary widely internationally.
- 2683 • Feedback from European computer manufacturers indicated that this facility was not
2684 consistent with European needs and that specification of such a facility was not a
2685 requirement for portability.

2686 **A.11.1 Interface Characteristics**2687 *A.11.1.1 Opening a Terminal Device File*

2688 There is no additional rationale provided for this section.

2689 *A.11.1.2 Process Groups*

2690 There is a potential race when the members of the foreground process group on a terminal leave
2691 that process group, either by exit or by changing process groups. After the last process exits the
2692 process group, but before the foreground process group ID of the terminal is changed (usually
2693 by a job control shell), it would be possible for a new process to be created with its process ID
2694 equal to the terminal's foreground process group ID. That process might then become the
2695 process group leader and accidentally be placed into the foreground on a terminal that was not
2696 necessarily its controlling terminal. As a result of this problem, the controlling terminal is
2697 defined to not have a foreground process group during this time.

2698 The cases where a controlling terminal has no foreground process group occur when all
2699 processes in the foreground process group either terminate and are waited for or join other
2700 process groups via *setpgid()* or *setsid()*. If the process group leader terminates, this is the first
2701 case described; if it leaves the process group via *setpgid()*, this is the second case described (a
2702 process group leader cannot successfully call *setsid()*). When one of those cases causes a
2703 controlling terminal to have no foreground process group, it has two visible effects on
2704 applications. The first is the value returned by *tcgetpgrp()*. The second (which occurs only in the
2705 case where the process group leader terminates) is the sending of signals in response to special
2706 input characters. The intent of IEEE Std 1003.1-2001 is that no process group be wrongly
2707 identified as the foreground process group by *tcgetpgrp()* or unintentionally receive signals
2708 because of placement into the foreground.

2709 In 4.3 BSD, the old process group ID continues to be used to identify the foreground process
2710 group and is returned by the function equivalent to *tcgetpgrp()*. In that implementation it is
2711 possible for a newly created process to be assigned the same value as a process ID and then form
2712 a new process group with the same value as a process group ID. The result is that the new
2713 process group would receive signals from this terminal for no apparent reason, and
2714 IEEE Std 1003.1-2001 precludes this by forbidding a process group from entering the foreground
2715 in this way. It would be more direct to place part of the requirement made by the last sentence
2716 under *fork()*, but there is no convenient way for that section to refer to the value that *tcgetpgrp()*
2717 returns, since in this case there is no process group and thus no process group ID.

2718 One possibility for a conforming implementation is to behave similarly to 4.3 BSD, but to
2719 prevent this reuse of the ID, probably in the implementation of *fork()*, as long as it is in use by
2720 the terminal.

2721 Another possibility is to recognize when the last process stops using the terminal's foreground
2722 process group ID, which is when the process group lifetime ends, and to change the terminal's
2723 foreground process group ID to a reserved value that is never used as a process ID or process
2724 group ID. (See the definition of *process group lifetime* in the definitions section.) The process ID
2725 can then be reserved until the terminal has another foreground process group.

2726 The 4.3 BSD implementation permits the leader (and only member) of the foreground process
2727 group to leave the process group by calling the equivalent of *setpgid()* and to later return,
2728 expecting to return to the foreground. There are no known application needs for this behavior,
2729 and IEEE Std 1003.1-2001 neither requires nor forbids it (except that it is forbidden for session
2730 leaders) by leaving it unspecified.

2731 A.11.1.3 *The Controlling Terminal*

2732 IEEE Std 1003.1-2001 does not specify a mechanism by which to allocate a controlling terminal.
2733 This is normally done by a system utility (such as *getty*) and is considered an administrative
2734 feature outside the scope of IEEE Std 1003.1-2001.

2735 Historical implementations allocate controlling terminals on certain *open()* calls. Since *open()* is
2736 part of POSIX.1, its behavior had to be dealt with. The traditional behavior is not required
2737 because it is not very straightforward or flexible for either implementations or applications.
2738 However, because of its prevalence, it was not practical to disallow this behavior either. Thus, a
2739 mechanism was standardized to ensure portable, predictable behavior in *open()*.

2740 Some historical implementations deallocate a controlling terminal on the last system-wide close.
2741 This behavior is neither required nor prohibited. Even on implementations that do provide this
2742 behavior, applications generally cannot depend on it due to its system-wide nature.

2743 A.11.1.4 *Terminal Access Control*

2744 The access controls described in this section apply only to a process that is accessing its
2745 controlling terminal. A process accessing a terminal that is not its controlling terminal is
2746 effectively treated the same as a member of the foreground process group. While this may seem
2747 unintuitive, note that these controls are for the purpose of job control, not security, and job
2748 control relates only to a process' controlling terminal. Normal file access permissions handle
2749 security.

2750 If the process calling *read()* or *write()* is in a background process group that is orphaned, it is not
2751 desirable to stop the process group, as it is no longer under the control of a job control shell that
2752 could put it into the foreground again. Accordingly, calls to *read()* or *write()* functions by such
2753 processes receive an immediate error return. This is different from 4.2 BSD, which kills orphaned
2754 processes that receive terminal stop signals.

2755 The foreground/background/orphaned process group check performed by the terminal driver
2756 must be repeatedly performed until the calling process moves into the foreground or until the
2757 process group of the calling process becomes orphaned. That is, when the terminal driver
2758 determines that the calling process is in the background and should receive a job control signal,
2759 it sends the appropriate signal (SIGTTIN or SIGTTOU) to every process in the process group of
2760 the calling process and then it allows the calling process to immediately receive the signal. The
2761 latter is typically performed by blocking the process so that the signal is immediately noticed.
2762 Note, however, that after the process finishes receiving the signal and control is returned to the
2763 driver, the terminal driver must re-execute the foreground/background/orphaned process
2764 group check. The process may still be in the background, either because it was continued in the
2765 background by a job control shell, or because it caught the signal and did nothing.

2766 The terminal driver repeatedly performs the foreground/background/orphaned process group
2767 checks whenever a process is about to access the terminal. In the case of *write()* or the control
2768 *tc*()* functions, the check is performed at the entry of the function. In the case of *read()*, the check
2769 is performed not only at the entry of the function, but also after blocking the process to wait for
2770 input characters (if necessary). That is, once the driver has determined that the process calling
2771 the *read()* function is in the foreground, it attempts to retrieve characters from the input queue. If
2772 the queue is empty, it blocks the process waiting for characters. When characters are available
2773 and control is returned to the driver, the terminal driver must return to the repeated
2774 foreground/background/orphaned process group check again. The process may have moved
2775 from the foreground to the background while it was blocked waiting for input characters.

2776 A.11.1.5 *Input Processing and Reading Data*

2777 There is no additional rationale provided for this section.

2778 A.11.1.6 *Canonical Mode Input Processing*

2779 The term “character” is intended here. ERASE should erase the last character, not the last byte.
2780 In the case of multi-byte characters, these two may be different.

2781 4.3 BSD has a WERASE character that erases the last “word” typed (but not any preceding
2782 <blank>s or <tab>s). A word is defined as a sequence of non-<blank>s, with <tab>s counted as
2783 <blank>s. Like ERASE, WERASE does not erase beyond the beginning of the line. This
2784 WERASE feature has not been specified in POSIX.1 because it is difficult to define in the
2785 international environment. It is only useful for languages where words are delimited by
2786 <blank>s. In some ideographic languages, such as Japanese and Chinese, words are not
2787 delimited at all. The WERASE character should presumably go back to the beginning of a
2788 sentence in those cases; practically, this means it would not be used much for those languages.

2789 It should be noted that there is a possible inherent deadlock if the application and
2790 implementation conflict on the value of {MAX_CANON}. With ICANON set (if IXOFF is
2791 enabled) and more than {MAX_CANON} characters transmitted without a <linefeed>,
2792 transmission will be stopped, the <linefeed> (or <carriage-return> when ICRLF is set) will never
2793 arrive, and the *read()* will never be satisfied.

2794 An application should not set IXOFF if it is using canonical mode unless it knows that (even in
2795 the face of a transmission error) the conditions described previously cannot be met or unless it is
2796 prepared to deal with the possible deadlock in some other way, such as timeouts.

2797 It should also be noted that this can be made to happen in non-canonical mode if the trigger
2798 value for sending IXOFF is less than VMIN and VTIME is zero.

2799 A.11.1.7 *Non-Canonical Mode Input Processing*

2800 Some points to note about MIN and TIME:

- 2801 1. The interactions of MIN and TIME are not symmetric. For example, when MIN>0 and
2802 TIME=0, TIME has no effect. However, in the opposite case where MIN=0 and TIME>0,
2803 both MIN and TIME play a role in that MIN is satisfied with the receipt of a single
2804 character.
- 2805 2. Also note that in case A (MIN>0, TIME>0), TIME represents an inter-character timer, while
2806 in case C (MIN=0, TIME>0), TIME represents a read timer.

2807 These two points highlight the dual purpose of the MIN/TIME feature. Cases A and B, where
2808 MIN>0, exist to handle burst-mode activity (for example, file transfer programs) where a
2809 program would like to process at least MIN characters at a time. In case A, the inter-character
2810 timer is activated by a user as a safety measure; in case B, it is turned off.

2811 Cases C and D exist to handle single-character timed transfers. These cases are readily adaptable
2812 to screen-based applications that need to know if a character is present in the input queue before
2813 refreshing the screen. In case C, the read is timed; in case D, it is not.

2814 Another important note is that MIN is always just a minimum. It does not denote a record
2815 length. That is, if a program does a read of 20 bytes, MIN is 10, and 25 characters are present, 20
2816 characters are returned to the user. In the special case of MIN=0, this still applies: if more than
2817 one character is available, they all will be returned immediately.

2818 **A.11.1.8 Writing Data and Output Processing**

2819 There is no additional rationale provided for this section.

2820 **A.11.1.9 Special Characters**

2821 There is no additional rationale provided for this section.

2822 **A.11.1.10 Modem Disconnect**

2823 There is no additional rationale provided for this section.

2824 **A.11.1.11 Closing a Terminal Device File**

2825 IEEE Std 1003.1-2001 does not specify that a *close()* on a terminal device file include the
2826 equivalent of a call to *tcfow(fd,TCOON)*.

2827 An implementation that discards output at the time *close()* is called after reporting the return
2828 value to the *write()* call that data was written does not conform with IEEE Std 1003.1-2001. An
2829 application has functions such as *tcdrain()*, *tcfush()*, and *tcfow()* available to obtain the detailed
2830 behavior it requires with respect to flushing of output.

2831 At the time of the last close on a terminal device, an application relinquishes any ability to exert
2832 flow control via *tcfow()*.

2833 **A.11.2 Parameters that Can be Set**2834 **A.11.2.1 The termios Structure**

2835 This structure is part of an interface that, in general, retains the historic grouping of flags.
2836 Although a more optimal structure for implementations may be possible, the degree of change
2837 to applications would be significantly larger.

2838 **A.11.2.2 Input Modes**

2839 Some historical implementations treated a long break as multiple events, as many as one per
2840 character time. The wording in POSIX.1 explicitly prohibits this.

2841 Although the *ISTRIP* flag is normally superfluous with today's terminal hardware and software,
2842 it is historically supported. Therefore, applications may be using *ISTRIP*, and there is no
2843 technical problem with supporting this flag. Also, applications may wish to receive only 7-bit
2844 input bytes and may not be connected directly to the hardware terminal device (for example,
2845 when a connection traverses a network).

2846 Also, there is no requirement in general that the terminal device ensures that high-order bits
2847 beyond the specified character size are cleared. *ISTRIP* provides this function for 7-bit
2848 characters, which are common.

2849 In dealing with multi-byte characters, the consequences of a parity error in such a character, or in
2850 an escape sequence affecting the current character set, are beyond the scope of POSIX.1 and are
2851 best dealt with by the application processing the multi-byte characters.

2852 A.11.2.3 *Output Modes*

2853 POSIX.1 does not describe postprocessing of output to a terminal or detailed control of that from
2854 a conforming application. (That is, translation of <newline> to <carriage-return> followed by
2855 <linefeed> or <tab> processing.) There is nothing that a conforming application should do to its
2856 output for a terminal because that would require knowledge of the operation of the terminal. It
2857 is the responsibility of the operating system to provide postprocessing appropriate to the output
2858 device, whether it is a terminal or some other type of device.

2859 Extensions to POSIX.1 to control the type of postprocessing already exist and are expected to
2860 continue into the future. The control of these features is primarily to adjust the interface between
2861 the system and the terminal device so the output appears on the display correctly. This should
2862 be set up before use by any application.

2863 In general, both the input and output modes should not be set absolutely, but rather modified
2864 from the inherited state.

2865 A.11.2.4 *Control Modes*

2866 This section could be misread that the symbol “CSIZE” is a title in the **termios** *c_flag* field.
2867 Although it does serve that function, it is also a required symbol, as a literal reading of POSIX.1
2868 (and the caveats about typography) would indicate.

2869 A.11.2.5 *Local Modes*

2870 Non-canonical mode is provided to allow fast bursts of input to be read efficiently while still
2871 allowing single-character input.

2872 The ECHONL function historically has been in many implementations. Since there seems to be
2873 no technical problem with supporting ECHONL, it is included in POSIX.1 to increase consensus.

2874 The alternate behavior possible when ECHOK or ECHOE are specified with ICANON is
2875 permitted as a compromise depending on what the actual terminal hardware can do. Erasing
2876 characters and lines is preferred, but is not always possible.

2877 A.11.2.6 *Special Control Characters*

2878 Permitting VMIN and VTIME to overlap with VEOF and VEOL was a compromise for historical
2879 implementations. Only when backwards-compatibility of object code is a serious concern to an
2880 implementor should an implementation continue this practice. Correct applications that work
2881 with the overlap (at the source level) should also work if it is not present, but not the reverse.

2882 **A.12 Utility Conventions**2883 **A.12.1 Utility Argument Syntax**

2884 The standard developers considered that recent trends toward diluting the SYNOPSIS sections
2885 of historical reference pages to the equivalent of:

2886 `command [options][operands]`

2887 were a disservice to the reader. Therefore, considerable effort was placed into rigorous
2888 definitions of all the command line arguments and their interrelationships. The relationships
2889 depicted in the synopses are normative parts of IEEE Std 1003.1-2001; this information is
2890 sometimes repeated in textual form, but that is only for clarity within context.

2891 The use of “undefined” for conflicting argument usage and for repeated usage of the same
2892 option is meant to prevent conforming applications from using conflicting arguments or
2893 repeated options unless specifically allowed (as is the case with *ls*, which allows simultaneous,
2894 repeated use of the *-C*, *-l*, and *-1* options). Many historical implementations will tolerate this
2895 usage, choosing either the first or the last applicable argument. This tolerance can continue, but
2896 conforming applications cannot rely upon it. (Other implementations may choose to print usage
2897 messages instead.)

2898 The use of “undefined” for conflicting argument usage also allows an implementation to make
2899 reasonable extensions to utilities where the implementor considers mutually-exclusive options
2900 according to IEEE Std 1003.1-2001 to have a sensible meaning and result.

2901 IEEE Std 1003.1-2001 does not define the result of a command when an option-argument or
2902 operand is not followed by ellipses and the application specifies more than one of that option-
2903 argument or operand. This allows an implementation to define valid (although non-standard)
2904 behavior for the utility when more than one such option or operand is specified.

2905 The following table summarizes the requirements for option-arguments:

	SYNOPSIS Shows:		
	<i>-a arg</i>	<i>-barg</i>	<i>-c[arg]</i>
Conforming application uses:	<i>-a arg</i>	<i>-barg</i>	<i>-carg</i> or <i>-c</i>
System supports:	<i>-a arg</i> and <i>-aarg</i>	<i>-b arg</i> and <i>-barg</i>	<i>-carg</i> and <i>-c</i>
Non-conforming applications may use:	<i>-aarg</i>	<i>-b arg</i>	N/A

2913 Allowing <blank>s after an option (that is, placing an option and its option-argument into
2914 separate argument strings) when IEEE Std 1003.1-2001 does not require it encourages portability
2915 of users, while still preserving backwards-compatibility of scripts. Inserting <blank>s between
2916 the option and the option-argument is preferred; however, historical usage has not been
2917 consistent in this area; therefore, <blank>s are required to be handled by all implementations,
2918 but implementations are also allowed to handle the historical syntax. Another justification for
2919 selecting the multiple-argument method was that the single-argument case is inherently
2920 ambiguous when the option-argument can legitimately be a null string.

2921 IEEE Std 1003.1-2001 explicitly states that digits are permitted as operands and option-
2922 arguments. The lower and upper bounds for the values of the numbers used for operands and
2923 option-arguments were derived from the ISO C standard values for {LONG_MIN} and
2924 {LONG_MAX}. The requirement on the standard utilities is that numbers in the specified range
2925 do not cause a syntax error, although the specification of a number need not be semantically
2926 correct for a particular operand or option-argument of a utility. For example, the specification of:

2927 dd obs=3000000000

2928 would yield undefined behavior for the application and could be a syntax error because the
2929 number 3 000 000 000 is outside of the range $-2\ 147\ 483\ 647$ to $+2\ 147\ 483\ 647$. On the other hand:

2930 dd obs=2000000000

2931 may cause some error, such as “blocksize too large”, rather than a syntax error.

2932 **A.12.2 Utility Syntax Guidelines**

2933 This section is based on the rules listed in the SVID. It was included for two reasons:

- 2934 1. The individual utility descriptions in the Shell and Utilities volume of
2935 IEEE Std 1003.1-2001, Chapter 4, Utilities needed a set of common (although not universal)
2936 actions on which they could anchor their descriptions of option and operand syntax. Most
2937 of the standard utilities actually do use these guidelines, and many of their historical
2938 implementations use the *getopt()* function for their parsing. Therefore, it was simpler to
2939 cite the rules and merely identify exceptions.
- 2940 2. Writers of conforming applications need suggested guidelines if the POSIX community is
2941 to avoid the chaos of historical UNIX system command syntax.

2942 It is recommended that all *future* utilities and applications use these guidelines to enhance “user
2943 portability”. The fact that some historical utilities could not be changed (to avoid breaking
2944 historical applications) should not deter this future goal.

2945 The voluntary nature of the guidelines is highlighted by repeated uses of the word *should*
2946 throughout. This usage should not be misinterpreted to imply that utilities that claim
2947 conformance in their OPTIONS sections do not always conform.

2948 Guidelines 1 and 2 are offered as guidance for locales using Latin alphabets. No
2949 recommendations are made by IEEE Std 1003.1-2001 concerning utility naming in other locales.

2950 In the Shell and Utilities volume of IEEE Std 1003.1-2001, Section 2.9.1, Simple Commands, it is
2951 further stated that a command used in the Shell Command Language cannot be named with a
2952 trailing colon.

2953 Guideline 3 was changed to allow alphanumeric characters (letters and digits) from the character
2954 set to allow compatibility with historical usage. Historical practice allows the use of digits
2955 wherever practical, and there are no portability issues that would prohibit the use of digits. In
2956 fact, from an internationalization viewpoint, digits (being non-language-dependent) are
2957 preferable over letters (a -2 is intuitively self-explanatory to any user, while in the $-f$ *filename*
2958 the letter ‘*f*’ is a mnemonic aid only to speakers of Latin-based languages where “filename”
2959 happens to translate to a word that begins with ‘*f*’). Since guideline 3 still retains the word
2960 “single”, multi-digit options are not allowed. Instances of historical utilities that used them have
2961 been marked obsolescent, with the numbers being changed from option names to option-
2962 arguments.

2963 It was difficult to achieve a satisfactory solution to the problem of name space in option
2964 characters. When the standard developers desired to extend the historical *cc* utility to accept
2965 ISO C standard programs, they found that all of the portable alphabet was already in use by
2966 various vendors. Thus, they had to devise a new name, *c89* (now superseded by *c99*), rather than
2967 something like *cc -X*. There were suggestions that implementors be restricted to providing
2968 extensions through various means (such as using a plus sign as the option delimiter or using
2969 option characters outside the alphanumeric set) that would reserve all of the remaining
2970 alphanumeric characters for future POSIX standards. These approaches were resisted because
2971 they lacked the historical style of UNIX systems. Furthermore, if a vendor-provided option

2972 should become commonly used in the industry, it would be a candidate for standardization. It
 2973 would be desirable to standardize such a feature using historical practice for the syntax (the
 2974 semantics can be standardized with any syntax). This would not be possible if the syntax was
 2975 one reserved for the vendor. However, since the standardization process may lead to minor
 2976 changes in the semantics, it may prove to be better for a vendor to use a syntax that will not be
 2977 affected by standardization.

2978 Guideline 8 includes the concept of comma-separated lists in a single argument. It is up to the
 2979 utility to parse such a list itself because *getopt()* just returns the single string. This situation was
 2980 retained so that certain historical utilities would not violate the guidelines. Applications
 2981 preparing for international use should be aware of an occasional problem with comma-
 2982 separated lists: in some locales, the comma is used as the radix character. Thus, if an application
 2983 is preparing operands for a utility that expects a comma-separated list, it should avoid
 2984 generating non-integer values through one of the means that is influenced by setting the
 2985 *LC_NUMERIC* variable (such as *awk*, *bc*, *printf*, or *printf()*).

2986 Applications calling any utility with a first operand starting with '-' should usually specify --,
 2987 as indicated by Guideline 10, to mark the end of the options. This is true even if the SYNOPSIS
 2988 in the Shell and Utilities volume of IEEE Std 1003.1-2001 does not specify any options;
 2989 implementations may provide options as extensions to the Shell and Utilities volume of
 2990 IEEE Std 1003.1-2001. The standard utilities that do not support Guideline 10 indicate that fact in
 2991 the OPTIONS section of the utility description.

2992 Guideline 11 was modified to clarify that the order of different options should not matter
 2993 relative to one another. However, the order of repeated options that also have option-arguments
 2994 may be significant; therefore, such options are required to be interpreted in the order that they
 2995 are specified. The *make* utility is an instance of a historical utility that uses repeated options in
 2996 which the order is significant. Multiple files are specified by giving multiple instances of the -f
 2997 option; for example:

```
2998     make -f common_header -f specific_rules target
```

2999 Guideline 13 does not imply that all of the standard utilities automatically accept the operand
 3000 '-' to mean standard input or output, nor does it specify the actions of the utility upon
 3001 encountering multiple '-' operands. It simply says that, by default, '-' operands are not used
 3002 for other purposes in the file reading or writing (but not when using *stat()*, *unlink()*, *touch*, and
 3003 so on) utilities. All information concerning actual treatment of the '-' operand is found in the
 3004 individual utility sections.

3005 An area of concern was that as implementations mature, implementation-defined utilities and
 3006 implementation-defined utility options will result. The idea was expressed that there needed to
 3007 be a standard way, say an environment variable or some such mechanism, to identify
 3008 implementation-defined utilities separately from standard utilities that may have the same
 3009 name. It was decided that there already exist several ways of dealing with this situation and that
 3010 it is outside of the scope to attempt to standardize in the area of non-standard items. A method
 3011 that exists on some historical implementations is the use of the so-called */local/bin* or
 3012 */usr/local/bin* directory to separate local or additional copies or versions of utilities. Another
 3013 method that is also used is to isolate utilities into completely separate domains. Still another
 3014 method to ensure that the desired utility is being used is to request the utility by its full
 3015 pathname. There are many approaches to this situation; the examples given above serve to
 3016 illustrate that there is more than one.

3017 A.13 Headers**3018 A.13.1 Format of Entries**

3019 Each header reference page has a common layout of sections describing the interface. This layout
3020 is similar to the manual page or “man” page format shipped with most UNIX systems, and each
3021 header has sections describing the SYNOPSIS and DESCRIPTION. These are the two sections
3022 that relate to conformance.

3023 Additional sections are informative, and add considerable information for the application
3024 developer. APPLICATION USAGE sections provide additional caveats, issues, and
3025 recommendations to the developer. RATIONALE sections give additional information on the
3026 decisions made in defining the interface.

3027 FUTURE DIRECTIONS sections act as pointers to related work that may impact the interface in
3028 the future, and often cautions the developer to architect the code to account for a change in this
3029 area. Note that a future directions statement should not be taken as a commitment to adopt a
3030 feature or interface in the future.

3031 The CHANGE HISTORY section describes when the interface was introduced, and how it has
3032 changed.

3033 Option labels and margin markings in the page can be useful in guiding the application
3034 developer.

3035 / *Rationale (Informative)*

3036 **Part B:**

3037 **System Interfaces**

3038 *The Open Group*
3039 *The Institute of Electrical and Electronics Engineers, Inc.*

Rationale for System Interfaces

3040

3041 **B.1 Introduction**

3042 **B.1.1 Scope**

3043 Refer to Section A.1.1 (on page 3).

3044 **B.1.2 Conformance**

3045 Refer to Section A.2 (on page 9).

3046 **B.1.3 Normative References**

3047 There is no additional rationale provided for this section.

3048 **B.1.4 Change History**

3049 The change history is provided as an informative section, to track changes from previous issues
3050 of IEEE Std 1003.1-2001.

3051 The following sections describe changes made to the System Interfaces volume of
3052 IEEE Std 1003.1-2001 since Issue 5 of the base document. The CHANGE HISTORY section for
3053 each entry details the technical changes that have been made to that entry from Issue 5. Changes
3054 between earlier issues of the base document and Issue 5 are not included.

3055 The change history between Issue 5 and Issue 6 also lists the changes since the
3056 ISO POSIX-1: 1996 standard.

3057 **Changes from Issue 5 to Issue 6 (IEEE Std 1003.1-2001)**

3058 The following list summarizes the major changes that were made in the System Interfaces
3059 volume of IEEE Std 1003.1-2001 from Issue 5 to Issue 6:

- 3060 • This volume of IEEE Std 1003.1-2001 is extensively revised so that it can be both an IEEE
3061 POSIX Standard and an Open Group Technical Standard.
- 3062 • The POSIX System Interfaces requirements incorporate support of FIPS 151-2.
- 3063 • The POSIX System Interfaces requirements are updated to align with some features of the
3064 Single UNIX Specification.
- 3065 • A RATIONALE section is added to each reference page.
- 3066 • Networking interfaces from the XNS, Issue 5.2 specification are incorporated.
- 3067 • IEEE Std 1003.1d-1999 is incorporated.
- 3068 • IEEE Std 1003.1j-2000 is incorporated.
- 3069 • IEEE Std 1003.1q-2000 is incorporated.
- 3070 • IEEE P1003.1a draft standard is incorporated.

- 3071 • Existing functionality is aligned with the ISO/IEC 9899: 1999 standard.
- 3072 • New functionality from the ISO/IEC 9899: 1999 standard is incorporated.
- 3073 • IEEE PASC Interpretations are applied.
- 3074 • The Open Group corrigenda and resolutions are applied.

3075 **New Features in Issue 6**

3076 The functions first introduced in Issue 6 (over the Issue 5 Base document) are listed in the table
3077 below:

3078
3079

3080
3081
3082
3083
3084
3085
3086
3087
3088
3089
3090
3091
3092
3093
3094
3095
3096
3097
3098
3099
3100
3101
3102
3103
3104
3105
3106
3107
3108
3109
3110
3111
3112
3113
3114
3115
3116

New Functions in Issue 6		
<i>acosf()</i>	<i>catanh()</i>	<i>cprojf()</i>
<i>acoshf()</i>	<i>catanl()</i>	<i>cprojl()</i>
<i>acoshl()</i>	<i>cbrtf()</i>	<i>creal()</i>
<i>acosl()</i>	<i>cbrtl()</i>	<i>crealf()</i>
<i>asinf()</i>	<i>ccos()</i>	<i>creall()</i>
<i>asinhf()</i>	<i>ccosf()</i>	<i>csin()</i>
<i>asinhl()</i>	<i>ccosh()</i>	<i>csinf()</i>
<i>asinl()</i>	<i>ccoshf()</i>	<i>csinh()</i>
<i>atan2f()</i>	<i>ccoshl()</i>	<i>csinhf()</i>
<i>atan2l()</i>	<i>ccosl()</i>	<i>csinhl()</i>
<i>atanf()</i>	<i>ceilf()</i>	<i>csinl()</i>
<i>atanhf()</i>	<i>ceil()</i>	<i>csqrt()</i>
<i>atanhl()</i>	<i>cexp()</i>	<i>csqrtf()</i>
<i>atanl()</i>	<i>cexpf()</i>	<i>csqrtl()</i>
<i>atoll()</i>	<i>cexpl()</i>	<i>ctan()</i>
<i>cabs()</i>	<i>cimag()</i>	<i>ctanf()</i>
<i>cabsf()</i>	<i>cimagf()</i>	<i>ctanh()</i>
<i>cabsl()</i>	<i>cimagl()</i>	<i>ctanhf()</i>
<i>cacos()</i>	<i>clock_getcpuclockid()</i>	<i>ctanhl()</i>
<i>cacosf()</i>	<i>clock_nanosleep()</i>	<i>ctanl()</i>
<i>cacosh()</i>	<i>clog()</i>	<i>erfcf()</i>
<i>cacoshf()</i>	<i>clogf()</i>	<i>erfcl()</i>
<i>cacoshl()</i>	<i>clogl()</i>	<i>erff()</i>
<i>cacosl()</i>	<i>conj()</i>	<i>erfl()</i>
<i>carg()</i>	<i>conjf()</i>	<i>exp2()</i>
<i>cargf()</i>	<i>conjl()</i>	<i>exp2f()</i>
<i>cargl()</i>	<i>copysign()</i>	<i>exp2l()</i>
<i>casin()</i>	<i>copysignf()</i>	<i>expf()</i>
<i>casinf()</i>	<i>copysignl()</i>	<i>expl()</i>
<i>casinh()</i>	<i>cosf()</i>	<i>expm1f()</i>
<i>casinhf()</i>	<i>coshf()</i>	<i>expm1l()</i>
<i>casinhl()</i>	<i>coshl()</i>	<i>fabsf()</i>
<i>casinl()</i>	<i>cosl()</i>	<i>fabsl()</i>
<i>catan()</i>	<i>cpow()</i>	<i>fdim()</i>
<i>catanf()</i>	<i>cpowf()</i>	<i>fdimf()</i>
<i>catanh()</i>	<i>cpowl()</i>	<i>fdiml()</i>
<i>catanhf()</i>	<i>cproj()</i>	<i>feclearexcept()</i>

3117
3118
3119
3120
3121
3122
3123
3124
3125
3126
3127
3128
3129
3130
3131
3132
3133
3134
3135
3136
3137
3138
3139
3140
3141
3142
3143
3144
3145
3146
3147
3148
3149
3150
3151
3152
3153
3154
3155
3156
3157
3158
3159
3160
3161
3162

New Functions in Issue 6

<i>fegetenv()</i>	<i>ldexpl()</i>	<i>posix_fallocate()</i>
<i>fegetexceptflag()</i>	<i>lgammaf()</i>	<i>posix_madvise()</i>
<i>fegetround()</i>	<i>lgammal()</i>	<i>posix_mem_offset()</i>
<i>fehldexcept()</i>	<i>llabs()</i>	<i>posix_memalign()</i>
<i>feraiseexcept()</i>	<i>lldiv()</i>	<i>posix_openpt()</i>
<i>fesetenv()</i>	<i>llrint()</i>	<i>posix_spawn()</i>
<i>fesetexceptflag()</i>	<i>llrintf()</i>	<i>posix_spawn_file_actions_addclose()</i>
<i>fesetround()</i>	<i>llrintl()</i>	<i>posix_spawn_file_actions_adddup2()</i>
<i>fetestexcept()</i>	<i>llround()</i>	<i>posix_spawn_file_actions_addopen()</i>
<i>feupdateenv()</i>	<i>llroundf()</i>	<i>posix_spawn_file_actions_destroy()</i>
<i>floorf()</i>	<i>llroundl()</i>	<i>posix_spawn_file_actions_init()</i>
<i>floorl()</i>	<i>log10f()</i>	<i>posix_spawnattr_destroy()</i>
<i>fna()</i>	<i>log10l()</i>	<i>posix_spawnattr_getflags()</i>
<i>fnaf()</i>	<i>log1pf()</i>	<i>posix_spawnattr_getpgroup()</i>
<i>fnal()</i>	<i>log1pl()</i>	<i>posix_spawnattr_getschedparam()</i>
<i>fnax()</i>	<i>log2()</i>	<i>posix_spawnattr_getschedpolicy()</i>
<i>fnaxf()</i>	<i>log2f()</i>	<i>posix_spawnattr_getsigdefault()</i>
<i>fnaxl()</i>	<i>log2l()</i>	<i>posix_spawnattr_getsigmask()</i>
<i>fnin()</i>	<i>logbf()</i>	<i>posix_spawnattr_init()</i>
<i>fninf()</i>	<i>logbl()</i>	<i>posix_spawnattr_setflags()</i>
<i>fninl()</i>	<i>logf()</i>	<i>posix_spawnattr_setpgroup()</i>
<i>fnodf()</i>	<i>logl()</i>	<i>posix_spawnattr_setschedparam()</i>
<i>fnodl()</i>	<i>lrint()</i>	<i>posix_spawnattr_setschedpolicy()</i>
<i>fpclassify()</i>	<i>lrintf()</i>	<i>posix_spawnattr_setsigdefault()</i>
<i>frexpf()</i>	<i>lrintl()</i>	<i>posix_spawnattr_setsigmask()</i>
<i>frexpl()</i>	<i>lround()</i>	<i>posix_spawnnp()</i>
<i>hypotf()</i>	<i>lroundf()</i>	<i>posix_trace_attr_destroy()</i>
<i>hypotl()</i>	<i>lroundl()</i>	<i>posix_trace_attr_getclockres()</i>
<i>ilogbf()</i>	<i>modff()</i>	<i>posix_trace_attr_getcreatetime()</i>
<i>ilogbl()</i>	<i>modfl()</i>	<i>posix_trace_attr_getgenversion()</i>
<i>imaxabs()</i>	<i>mq_timedreceive()</i>	<i>posix_trace_attr_getinherited()</i>
<i>imaxdiv()</i>	<i>mq_timedsend()</i>	<i>posix_trace_attr_getlogfullpolicy()</i>
<i>isblank()</i>	<i>nan()</i>	<i>posix_trace_attr_getlogsize()</i>
<i>isfinite()</i>	<i>nanf()</i>	<i>posix_trace_attr_getmaxdatasize()</i>
<i>isgreater()</i>	<i>nanl()</i>	<i>posix_trace_attr_getmaxsystemeventsz()</i>
<i>isgreaterequal()</i>	<i>nearbyint()</i>	<i>posix_trace_attr_getmaxuserereventsiz()</i>
<i>isinf()</i>	<i>nearbyintf()</i>	<i>posix_trace_attr_getname()</i>
<i>isless()</i>	<i>nearbyintl()</i>	<i>posix_trace_attr_getstreamfullpolicy()</i>
<i>islessequal()</i>	<i>nextafterf()</i>	<i>posix_trace_attr_getstreamsize()</i>
<i>islessgreater()</i>	<i>nextafterl()</i>	<i>posix_trace_attr_init()</i>
<i>isnormal()</i>	<i>nexttoward()</i>	<i>posix_trace_attr_setinherited()</i>
<i>isunordered()</i>	<i>nexttowardf()</i>	<i>posix_trace_attr_setlogfullpolicy()</i>
<i>iswblank()</i>	<i>nexttowardl()</i>	<i>posix_trace_attr_setlogsize()</i>
<i>ldexpf()</i>	<i>posix_fadvise()</i>	<i>posix_trace_create()</i>

3163
3164
3165
3166
3167
3168
3169
3170
3171
3172
3173
3174
3175
3176
3177
3178
3179
3180
3181
3182
3183
3184
3185
3186
3187
3188
3189
3190
3191
3192
3193
3194
3195
3196
3197
3198
3199
3200
3201
3202
3203

New Functions in Issue 6

<i>posix_trace_attr_setmaxdatasize()</i>	<i>pthread_barrier_destroy()</i>	<i>signbit()</i>
<i>posix_trace_attr_setname()</i>	<i>pthread_barrier_init()</i>	<i>sinf()</i>
<i>posix_trace_attr_setstreamfullpolicy()</i>	<i>pthread_barrier_wait()</i>	<i>sinhf()</i>
<i>posix_trace_attr_setstreamsize()</i>	<i>pthread_barrierattr_destroy()</i>	<i>sinhl()</i>
<i>posix_trace_clear()</i>	<i>pthread_barrierattr_getpshared()</i>	<i>sinl()</i>
<i>posix_trace_close()</i>	<i>pthread_barrierattr_init()</i>	<i>socketmark()</i>
<i>posix_trace_create_withlog()</i>	<i>pthread_barrierattr_setpshared()</i>	<i>sqrtf()</i>
<i>posix_trace_event()</i>	<i>pthread_condattr_getclock()</i>	<i>sqrtl()</i>
<i>posix_trace_eventid_equal()</i>	<i>pthread_condattr_setclock()</i>	<i>strerror_r()</i>
<i>posix_trace_eventid_get_name()</i>	<i>pthread_getcpuclockid()</i>	<i>strtoimax()</i>
<i>posix_trace_eventid_open()</i>	<i>pthread_mutex_timedlock()</i>	<i>strtoll()</i>
<i>posix_trace_eventset_add()</i>	<i>pthread_rwlock_timedrdlock()</i>	<i>strtoull()</i>
<i>posix_trace_eventset_del()</i>	<i>pthread_rwlock_timedwrlock()</i>	<i>strtoumax()</i>
<i>posix_trace_eventset_empty()</i>	<i>pthread_setschedprio()</i>	<i>tanf()</i>
<i>posix_trace_eventset_fill()</i>	<i>pthread_spin_destroy()</i>	<i>tanhf()</i>
<i>posix_trace_eventset_ismember()</i>	<i>pthread_spin_init()</i>	<i>tanhl()</i>
<i>posix_trace_eventtypelist_getnext_id()</i>	<i>pthread_spin_lock()</i>	<i>tanl()</i>
<i>posix_trace_eventtypelist_rewind()</i>	<i>pthread_spin_trylock()</i>	<i>tgamma()</i>
<i>posix_trace_flush()</i>	<i>pthread_spin_unlock()</i>	<i>tgammaf()</i>
<i>posix_trace_get_attr()</i>	<i>remainderf()</i>	<i>tgamma1()</i>
<i>posix_trace_get_filter()</i>	<i>remainderl()</i>	<i>trunc()</i>
<i>posix_trace_get_status()</i>	<i>remquo()</i>	<i>truncf()</i>
<i>posix_trace_getnext_event()</i>	<i>remquof()</i>	<i>truncl()</i>
<i>posix_trace_open()</i>	<i>remquol()</i>	<i>unsetenv()</i>
<i>posix_trace_rewind()</i>	<i>rintf()</i>	<i>vfprintf()</i>
<i>posix_trace_set_filter()</i>	<i>rintl()</i>	<i>vfprintf()</i>
<i>posix_trace_shutdown()</i>	<i>round()</i>	<i>vwscanf()</i>
<i>posix_trace_start()</i>	<i>roundf()</i>	<i>vprintf()</i>
<i>posix_trace_stop()</i>	<i>roundl()</i>	<i>vscanf()</i>
<i>posix_trace_timedgetnext_event()</i>	<i>scalbln()</i>	<i>vsprintf()</i>
<i>posix_trace_trid_eventid_open()</i>	<i>scalblnf()</i>	<i>vsscanf()</i>
<i>posix_trace_trygetnext_event()</i>	<i>scalblnl()</i>	<i>vswscanf()</i>
<i>posix_typed_mem_get_info()</i>	<i>scalbn()</i>	<i>vwscanf()</i>
<i>posix_typed_mem_open()</i>	<i>scalbnf()</i>	<i>wcstoimax()</i>
<i>powf()</i>	<i>scalbnl()</i>	<i>wcstoll()</i>
<i>powl()</i>	<i>sem_timedwait()</i>	<i>wcstoull()</i>
<i>pselect()</i>	<i>setgid()</i>	<i>wcstoumax()</i>
<i>pthread_attr_getstack()</i>	<i>setenv()</i>	
<i>pthread_attr_setstack()</i>	<i>setuid()</i>	

3204 The following new headers are introduced in Issue 6:

3205

3206

3207

3208

3209

New Headers in Issue 6		
<complex.h>	<spawn.h>	<tgmath.h>
<fenv.h>	<stdbool.h>	<trace.h>
<net/if.h>	<stdint.h>	

3210 The following table lists the functions and symbols from the XSI extension. These are new since
 3211 the ISO POSIX-1: 1996 standard.

3212
 3213

New XSI Functions and Symbols in Issue 6			
3214	<i>_longjmp()</i>	<i>getcontext()</i>	<i>msgget()</i>
3215	<i>_setjmp()</i>	<i>getdate()</i>	<i>msgrcv()</i>
3216	<i>_tolower()</i>	<i>getgrent()</i>	<i>msgsnd()</i>
3217	<i>_toupper()</i>	<i>gethostid()</i>	<i>nftw()</i>
3218	<i>a64l()</i>	<i>getitimer()</i>	<i>nice()</i>
3219	<i>basename()</i>	<i>getpgid()</i>	<i>nl_langinfo()</i>
3220	<i>bcmp()</i>	<i>getpmsg()</i>	<i>nrnd48()</i>
3221	<i>bcopy()</i>	<i>getpriority()</i>	<i>openlog()</i>
3222	<i>bzero()</i>	<i>getpwent()</i>	<i>poll()</i>
3223	<i>catclose()</i>	<i>getrlimit()</i>	<i>posix_openpt()</i>
3224	<i>catgets()</i>	<i>getrusage()</i>	<i>pread()</i>
3225	<i>catopen()</i>	<i>getsid()</i>	<i>pthread_attr_getguardsize()</i>
3226	<i>closelog()</i>	<i>getsubopt()</i>	<i>pthread_attr_setguardsize()</i>
3227	<i>crypt()</i>	<i>gettimeofday()</i>	<i>pthread_attr_setstack()</i>
3228	<i>daylight</i>	<i>getutxent()</i>	<i>pthread_getconcurrency()</i>
3229	<i>dbm_clearerr()</i>	<i>getutxid()</i>	<i>pthread_mutexattr_gettype()</i>
3230	<i>dbm_close()</i>	<i>getutxline()</i>	<i>pthread_mutexattr_settype()</i>
3231	<i>dbm_delete()</i>	<i>getwd()</i>	<i>pthread_rwlockattr_init()</i>
3232	<i>dbm_error()</i>	<i>grantpt()</i>	<i>pthread_rwlockattr_setpshared()</i>
3233	<i>dbm_fetch()</i>	<i>hcreate()</i>	<i>pthread_setconcurrency()</i>
3234	<i>dbm_firstkey()</i>	<i>hdestroy()</i>	<i>ptsname()</i>
3235	<i>dbm_nextkey()</i>	<i>hsearch()</i>	<i>putenv()</i>
3236	<i>dbm_open()</i>	<i>iconv()</i>	<i>pututxline()</i>
3237	<i>dbm_store()</i>	<i>iconv_close()</i>	<i>pwrite()</i>
3238	<i>dirname()</i>	<i>iconv_open()</i>	<i>random()</i>
3239	<i>dlclose()</i>	<i>index()</i>	<i>readv()</i>
3240	<i>dLError()</i>	<i>initstate()</i>	<i>realpath()</i>
3241	<i>dlopen()</i>	<i>insque()</i>	<i>remque()</i>
3242	<i>dlsym()</i>	<i>isascii()</i>	<i>rindex()</i>
3243	<i>drand48()</i>	<i>jrand48()</i>	<i>seed48()</i>
3244	<i>ecvt()</i>	<i>killpg()</i>	<i>seekdir()</i>
3245	<i>encrypt()</i>	<i>l64a()</i>	<i>semctl()</i>
3246	<i>endgrent()</i>	<i>lchown()</i>	<i>semget()</i>
3247	<i>endpwent()</i>	<i>lcong48()</i>	<i>semop()</i>
3248	<i>endutxent()</i>	<i>lfind()</i>	<i>setcontext()</i>
3249	<i>erand48()</i>	<i>lockf()</i>	<i>setgrent()</i>
3250	<i>fchdir()</i>	<i>lrnd48()</i>	<i>setitimer()</i>
3251	<i>fcvt()</i>	<i>lsearch()</i>	<i>setkey()</i>
3252	<i>ffs()</i>	<i>makecontext()</i>	<i>setlogmask()</i>
3253	<i>fntmsg()</i>	<i>memccpy()</i>	<i>setpgrp()</i>
3254	<i>fstatvfs()</i>	<i>mknod()</i>	<i>setpriority()</i>
3255	<i>ftime()</i>	<i>mkstemp()</i>	<i>setpwent()</i>
3256	<i>ftok()</i>	<i>mktemp()</i>	<i>setregid()</i>
3257	<i>ftw()</i>	<i>mrnd48()</i>	<i>setreuid()</i>
3258	<i>gcvt()</i>	<i>msgctl()</i>	<i>setrlimit()</i>
			<i>setstate()</i>
			<i>setutxent()</i>
			<i>shmat()</i>
			<i>shmctl()</i>
			<i>shmdt()</i>
			<i>shmget()</i>
			<i>sigaltstack()</i>
			<i>sighold()</i>
			<i>sigignore()</i>
			<i>siginterrupt()</i>
			<i>sigpause()</i>
			<i>sigrelse()</i>
			<i>sigset()</i>
			<i>srand48()</i>
			<i>srandom()</i>
			<i>statvfs()</i>
			<i>strcasecmp()</i>
			<i>strdup()</i>
			<i>strfmon()</i>
			<i>strncasecmp()</i>
			<i>strptime()</i>
			<i>swab()</i>
			<i>swapcontext()</i>
			<i>sync()</i>
			<i>syslog()</i>
			<i>tcgetsid()</i>
			<i>tdelete()</i>
			<i>telldir()</i>
			<i>tempnam()</i>
			<i>tfind()</i>
			<i>timezone</i>
			<i>toascii()</i>
			<i>truncate()</i>
			<i>tsearch()</i>
			<i>twalk()</i>
			<i>ulimit()</i>
			<i>unlockpt()</i>
			<i>utimes()</i>
			<i>waitid()</i>
			<i>wcswcs()</i>
			<i>wcswidth()</i>
			<i>wcwidth()</i>
			<i>writev()</i>

3259 The following table lists the headers from the XSI extension. These are new since the
3260 ISO POSIX-1:1996 standard.

3261
3262

New XSI Headers in Issue 6		
<cpio.h>	<poll.h>	<sys/statvfs.h>
<dlfcn.h>	<search.h>	<sys/time.h>
<fmtmsg.h>	<strings.h>	<sys/timeb.h>
<ftw.h>	<stropts.h>	<sys/uio.h>
<iconv.h>	<sys/ipc.h>	<syslog.h>
<langinfo.h>	<sys/mman.h>	<ucontext.h>
<libgen.h>	<sys/msg.h>	<ulimit.h>
<monetary.h>	<sys/resource.h>	<utmpx.h>
<ndbm.h>	<sys/sem.h>	
<nl_types.h>	<sys/shm.h>	

3273 **B.1.5 Terminology**

3274 Refer to Section A.1.4 (on page 5).

3275 **B.1.6 Definitions**

3276 Refer to Section A.3 (on page 13).

3277 **B.1.7 Relationship to Other Formal Standards**

3278 There is no additional rationale provided for this section.

3279 **B.1.8 Portability**

3280 Refer to Section A.1.5 (on page 8).

3281 *B.1.8.1 Codes*

3282 Refer to Section A.1.5.1 (on page 8).

3283 **B.1.9 Format of Entries**

3284 Each system interface reference page has a common layout of sections describing the interface.
3285 This layout is similar to the manual page or “man” page format shipped with most UNIX
3286 systems, and each header has sections describing the SYNOPSIS, DESCRIPTION, RETURN
3287 VALUE, and ERRORS. These are the four sections that relate to conformance.

3288 Additional sections are informative, and add considerable information for the application
3289 developer. EXAMPLES sections provide example usage. APPLICATION USAGE sections
3290 provide additional caveats, issues, and recommendations to the developer. RATIONALE
3291 sections give additional information on the decisions made in defining the interface.

3292 FUTURE DIRECTIONS sections act as pointers to related work that may impact the interface in
3293 the future, and often cautions the developer to architect the code to account for a change in this
3294 area. Note that a future directions statement should not be taken as a commitment to adopt a
3295 feature or interface in the future.

3296 The CHANGE HISTORY section describes when the interface was introduced, and how it has
3297 changed.

3298 Option labels and margin markings in the page can be useful in guiding the application
3299 developer.

3300 **B.2 General Information**

3301 **B.2.1 Use and Implementation of Functions**

3302 The information concerning the use of functions was adapted from a description in the ISO C
3303 standard. Here is an example of how an application program can protect itself from functions
3304 that may or may not be macros, rather than true functions:

3305 The *atoi()* function may be used in any of several ways:

- 3306 • By use of its associated header (possibly generating a macro expansion):

```
3307     #include <stdlib.h>
3308     /* ... */
3309     i = atoi(str);
```

- 3310 • By use of its associated header (assuredly generating a true function call):

```
3311     #include <stdlib.h>
3312     #undef atoi
3313     /* ... */
3314     i = atoi(str);
```

3315 or:

```
3316     #include <stdlib.h>
3317     /* ... */
3318     i = (atoi) (str);
```

- 3319 • By explicit declaration:

```
3320     extern int atoi (const char *);
3321     /* ... */
3322     i = atoi(str);
```

- 3323 • By implicit declaration:

```
3324     /* ... */
3325     i = atoi(str);
```

3326 (Assuming no function prototype is in scope. This is not allowed by the ISO C standard for
3327 functions with variable arguments; furthermore, parameter type conversion “widening” is
3328 subject to different rules in this case.)

3329 Note that the ISO C standard reserves names starting with ‘_’ for the compiler. Therefore, the
3330 compiler could, for example, implement an intrinsic, built-in function *_asm_builtin_atoi()*, which
3331 it recognized and expanded into inline assembly code. Then, in *<stdlib.h>*, there could be the
3332 following:

```
3333     #define atoi(X) _asm_builtin_atoi(X)
```

3334 The user’s “normal” call to *atoi()* would then be expanded inline, but the implementor would
3335 also be required to provide a callable function named *atoi()* for use when the application
3336 requires it; for example, if its address is to be stored in a function pointer variable.

3337 **B.2.2 The Compilation Environment**3338 *B.2.2.1 POSIX.1 Symbols*

3339 This and the following section address the issue of “name space pollution”. The ISO C standard
 3340 requires that the name space beyond what it reserves not be altered except by explicit action of
 3341 the application writer. This section defines the actions to add the POSIX.1 symbols for those
 3342 headers where both the ISO C standard and POSIX.1 need to define symbols, and also where the
 3343 XSI extension extends the base standard.

3344 When headers are used to provide symbols, there is a potential for introducing symbols that the
 3345 application writer cannot predict. Ideally, each header should only contain one set of symbols,
 3346 but this is not practical for historical reasons. Thus, the concept of feature test macros is
 3347 included. Two feature test macros are explicitly defined by IEEE Std 1003.1-2001; it is expected
 3348 that future revisions may add to this.

3349 **Note:** Feature test macros allow an application to announce to the implementation its desire to have
 3350 certain symbols and prototypes exposed. They should not be confused with the version test
 3351 macros and constants for options in `<unistd.h>` which are the implementation’s way of
 3352 announcing functionality to the application.

3353 It is further intended that these feature test macros apply only to the headers specified by
 3354 IEEE Std 1003.1-2001. Implementations are expressly permitted to make visible symbols not
 3355 specified by IEEE Std 1003.1-2001, within both POSIX.1 and other headers, under the control of
 3356 feature test macros that are not defined by IEEE Std 1003.1-2001.

3357 **The `_POSIX_C_SOURCE` Feature Test Macro**

3358 Since `_POSIX_SOURCE` specified by the POSIX.1-1990 standard did not have a value associated
 3359 with it, the `_POSIX_C_SOURCE` macro replaces it, allowing an application to inform the system
 3360 of the revision of the standard to which it conforms. This symbol will allow implementations to
 3361 support various revisions of IEEE Std 1003.1-2001 simultaneously. For instance, when either
 3362 `_POSIX_SOURCE` is defined or `_POSIX_C_SOURCE` is defined as 1, the system should make
 3363 visible the same name space as permitted and required by the POSIX.1-1990 standard. When
 3364 `_POSIX_C_SOURCE` is defined, the state of `_POSIX_SOURCE` is completely irrelevant.

3365 It is expected that C bindings to future POSIX standards will define new values for
 3366 `_POSIX_C_SOURCE`, with each new value reserving the name space for that new standard, plus
 3367 all earlier POSIX standards.

3368 **The `_XOPEN_SOURCE` Feature Test Macro**

3369 The feature test macro `_XOPEN_SOURCE` is provided as the announcement mechanism for the
 3370 application that it requires functionality from the Single UNIX Specification. `_XOPEN_SOURCE`
 3371 must be defined to the value 600 before the inclusion of any header to enable the functionality in
 3372 the Single UNIX Specification. Its definition subsumes the use of `_POSIX_SOURCE` and
 3373 `_POSIX_C_SOURCE`.

3374 An extract of code from a conforming application, that appears before any `#include` statements,
 3375 is given below:

```
3376 #define _XOPEN_SOURCE 600 /* Single UNIX Specification, Version 3 */
3377 #include ...
```

3378 Note that the definition of `_XOPEN_SOURCE` with the value 600 makes the definition of
 3379 `_POSIX_C_SOURCE` redundant and it can safely be omitted.

3380 B.2.2.2 The Name Space

3381 The reservation of identifiers is paraphrased from the ISO C standard. The text is included
 3382 because it needs to be part of IEEE Std 1003.1-2001, regardless of possible changes in future
 3383 versions of the ISO C standard.

3384 These identifiers may be used by implementations, particularly for feature test macros.
 3385 Implementations should not use feature test macro names that might be reasonably used by a
 3386 standard.

3387 Including headers more than once is a reasonably common practice, and it should be carried
 3388 forward from the ISO C standard. More significantly, having definitions in more than one
 3389 header is explicitly permitted. Where the potential declaration is “benign” (the same definition
 3390 twice) the declaration can be repeated, if that is permitted by the compiler. (This is usually true
 3391 of macros, for example.) In those situations where a repetition is not benign (for example,
 3392 **typedefs**), conditional compilation must be used. The situation actually occurs both within the
 3393 ISO C standard and within POSIX.1: **time_t** should be in **<sys/types.h>**, and the ISO C standard
 3394 mandates that it be in **<time.h>**.

3395 The area of name space pollution *versus* additions to structures is difficult because of the macro
 3396 structure of C. The following discussion summarizes all the various problems with and
 3397 objections to the issue.

3398 Note the phrase “user-defined macro”. Users are not permitted to define macro names (or any
 3399 other name) beginning with “_**[A-Z_]**”. Thus, the conflict cannot occur for symbols reserved
 3400 to the vendor’s name space, and the permission to add fields automatically applies, without
 3401 qualification, to those symbols.

3402 1. Data structures (and unions) need to be defined in headers by implementations to meet
 3403 certain requirements of POSIX.1 and the ISO C standard.

3404 2. The structures defined by POSIX.1 are typically minimal, and any practical
 3405 implementation would wish to add fields to these structures either to hold additional
 3406 related information or for backwards-compatibility (or both). Future standards (and *de*
 3407 *facto* standards) would also wish to add to these structures. Issues of field alignment make
 3408 it impractical (at least in the general case) to simply omit fields when they are not defined
 3409 by the particular standard involved.

3410 The **dirent** structure is an example of such a minimal structure (although one could argue
 3411 about whether the other fields need visible names). The **st_rdev** field of most
 3412 implementations’ **stat** structure is a common example where extension is needed and
 3413 where a conflict could occur.

3414 3. Fields in structures are in an independent name space, so the addition of such fields
 3415 presents no problem to the C language itself in that such names cannot interact with
 3416 identically named user symbols because access is qualified by the specific structure name.

3417 4. There is an exception to this: macro processing is done at a lexical level. Thus, symbols
 3418 added to a structure might be recognized as user-provided macro names at the location
 3419 where the structure is declared. This only can occur if the user-provided name is declared
 3420 as a macro before the header declaring the structure is included. The user’s use of the name
 3421 after the declaration cannot interfere with the structure because the symbol is hidden and
 3422 only accessible through access to the structure. Presumably, the user would not declare
 3423 such a macro if there was an intention to use that field name.

3424 5. Macros from the same or a related header might use the additional fields in the structure,
 3425 and those field names might also collide with user macros. Although this is a less frequent
 3426 occurrence, since macros are expanded at the point of use, no constraint on the order of use

3427 of names can apply.

3428 6. An “obvious” solution of using names in the reserved name space and then redefining
 3429 them as macros when they should be visible does not work because this has the effect of
 3430 exporting the symbol into the general name space. For example, given a (hypothetical)
 3431 system-provided header `<h.h>`, and two parts of a C program in `a.c` and `b.c`, in header
 3432 `<h.h>`:

```
3433     struct foo {
3434         int __i;
3435     }
3436
3437     #ifdef _FEATURE_TEST
3438     #define i __i;
3439     #endif
```

3439 In file `a.c`:

```
3440     #include h.h
3441     extern int i;
3442     ...
```

3443 In file `b.c`:

```
3444     extern int i;
3445     ...
```

3446 The symbol that the user thinks of as `i` in both files has an external name of `__i` in `a.c`; the
 3447 same symbol `i` in `b.c` has an external name `i` (ignoring any hidden manipulations the
 3448 compiler might perform on the names). This would cause a mysterious name resolution
 3449 problem when `a.o` and `b.o` are linked.

3450 Simply avoiding definition then causes alignment problems in the structure.

3451 A structure of the form:

```
3452     struct foo {
3453         union {
3454             int __i;
3455             #ifdef _FEATURE_TEST
3456             int i;
3457             #endif
3458         } __ii;
3459     }
```

3460 does not work because the name of the logical field `i` is `__ii.i`, and introduction of a macro
 3461 to restore the logical name immediately reintroduces the problem discussed previously
 3462 (although its manifestation might be more immediate because a syntax error would result
 3463 if a recursive macro did not cause it to fail first).

3464 7. A more workable solution would be to declare the structure:

```

3465     struct foo {
3466         #ifdef _FEATURE_TEST
3467             int i;
3468         #else
3469             int __i;
3470         #endif
3471     }

```

3472 However, if a macro (particularly one required by a standard) is to be defined that uses
 3473 this field, two must be defined: one that uses *i*, the other that uses *__i*. If more than one
 3474 additional field is used in a macro and they are conditional on distinct combinations of
 3475 features, the complexity goes up as 2^n .

3476 All this leaves a difficult situation: vendors must provide very complex headers to deal with
 3477 what is conceptually simple and safe—adding a field to a structure. It is the possibility of user-
 3478 provided macros with the same name that makes this difficult.

3479 Several alternatives were proposed that involved constraining the user's access to part of the
 3480 name space available to the user (as specified by the ISO C standard). In some cases, this was
 3481 only until all the headers had been included. There were two proposals discussed that failed to
 3482 achieve consensus:

- 3483 1. Limiting it for the whole program.
- 3484 2. Restricting the use of identifiers containing only uppercase letters until after all system
 3485 headers had been included. It was also pointed out that because macros might wish to
 3486 access fields of a structure (and macro expansion occurs totally at point of use) restricting
 3487 names in this way would not protect the macro expansion, and thus the solution was
 3488 inadequate.

3489 It was finally decided that reservation of symbols would occur, but as constrained.

3490 The current wording also allows the addition of fields to a structure, but requires that user
 3491 macros of the same name not interfere. This allows vendors to do one of the following:

- 3492 • Not create the situation (do not extend the structures with user-accessible names or use the
 3493 solution in (7) above)
- 3494 • Extend their compilers to allow some way of adding names to structures and macros safely

3495 There are at least two ways that the compiler might be extended: add new preprocessor
 3496 directives that turn off and on macro expansion for certain symbols (without changing the value
 3497 of the macro) and a function or lexical operation that suppresses expansion of a word. The latter
 3498 seems more flexible, particularly because it addresses the problem in macros as well as in
 3499 declarations.

3500 The following seems to be a possible implementation extension to the C language that will do
 3501 this: any token that during macro expansion is found to be preceded by three '#' symbols shall
 3502 not be further expanded in exactly the same way as described for macros that expand to their
 3503 own name as in Section 3.8.3.4 of the ISO C standard. A vendor may also wish to implement this
 3504 as an operation that is lexically a function, which might be implemented as:

```

3505     #define __safe_name(x) ###x

```

3506 Using a function notation would insulate vendors from changes in standards until such a
 3507 functionality is standardized (if ever). Standardization of such a function would be valuable
 3508 because it would then permit third parties to take advantage of it portably in software they may
 3509 supply.

3510 The symbols that are “explicitly permitted, but not required by IEEE Std 1003.1-2001” include
 3511 those classified below. (That is, the symbols classified below might, but are not required to, be
 3512 present when `_POSIX_C_SOURCE` is defined to have the value 200112L.)

3513 • Symbols in `<limits.h>` and `<unistd.h>` that are defined to indicate support for options or
 3514 limits that are constant at compile-time

3515 • Symbols in the name space reserved for the implementation by the ISO C standard

3516 • Symbols in a name space reserved for a particular type of extension (for example, type names
 3517 ending with `_t` in `<sys/types.h>`)

3518 • Additional members of structures or unions whose names do not reduce the name space
 3519 reserved for applications

3520 Since both implementations and future revisions of IEEE Std 1003.1 and other POSIX standards
 3521 may use symbols in the reserved spaces described in these tables, there is a potential for name
 3522 space clashes. To avoid future name space clashes when adding symbols, implementations
 3523 should not use the `posix_`, `POSIX_`, or `_POSIX_` prefixes.

3524 B.2.3 Error Numbers

3525 It was the consensus of the standard developers that to allow the conformance document to
 3526 state that an error occurs and under what conditions, but to disallow a statement that it never
 3527 occurs, does not make sense. It could be implied by the current wording that this is allowed, but
 3528 to reduce the possibility of future interpretation requests, it is better to make an explicit
 3529 statement.

3530 The ISO C standard requires that `errno` be an assignable lvalue. Originally, the definition in
 3531 POSIX.1 was stricter than that in the ISO C standard, `extern int errno`, in order to support
 3532 historical usage. In a multi-threaded environment, implementing `errno` as a global variable
 3533 results in non-deterministic results when accessed. It is required, however, that `errno` work as a
 3534 per-thread error reporting mechanism. In order to do this, a separate `errno` value has to be
 3535 maintained for each thread. The following section discusses the various alternative solutions
 3536 that were considered.

3537 In order to avoid this problem altogether for new functions, these functions avoid using `errno`
 3538 and, instead, return the error number directly as the function return value; a return value of zero
 3539 indicates that no error was detected.

3540 For any function that can return errors, the function return value is not used for any purpose
 3541 other than for reporting errors. Even when the output of the function is scalar, it is passed
 3542 through a function argument. While it might have been possible to allow some scalar outputs to
 3543 be coded as negative function return values and mixed in with positive error status returns, this
 3544 was rejected—using the return value for a mixed purpose was judged to be of limited use and
 3545 error prone.

3546 Checking the value of `errno` alone is not sufficient to determine the existence or type of an error,
 3547 since it is not required that a successful function call clear `errno`. The variable `errno` should only
 3548 be examined when the return value of a function indicates that the value of `errno` is meaningful.
 3549 In that case, the function is required to set the variable to something other than zero.

3550 The variable `errno` is never set to zero by any function call; to do so would contradict the ISO C
 3551 standard.

3552 POSIX.1 requires (in the ERRORS sections of function descriptions) certain error values to be set
 3553 in certain conditions because many existing applications depend on them. Some error numbers,
 3554 such as [EFAULT], are entirely implementation-defined and are noted as such in their

3555 description in the ERRORS section. This section otherwise allows wide latitude to the
 3556 implementation in handling error reporting.

3557 Some of the ERRORS sections in IEEE Std 1003.1-2001 have two subsections. The first:

3558 “The function shall fail if:”

3559 could be called the “mandatory” section.

3560 The second:

3561 “The function may fail if:”

3562 could be informally known as the “optional” section.

3563 Attempting to infer the quality of an implementation based on whether it detects optional error
 3564 conditions is not useful.

3565 Following each one-word symbolic name for an error, there is a description of the error. The
 3566 rationale for some of the symbolic names follows:

3567 [ECANCELED] This spelling was chosen as being more common.

3568 [EFAULT] Most historical implementations do not catch an error and set *errno* when an
 3569 invalid address is given to the functions *wait()*, *time()*, or *times()*. Some
 3570 implementations cannot reliably detect an invalid address. And most systems
 3571 that detect invalid addresses will do so only for a system call, not for a library
 3572 routine.

3573 [EFTYPE] This error code was proposed in earlier proposals as “Inappropriate operation
 3574 for file type”, meaning that the operation requested is not appropriate for the
 3575 file specified in the function call. This code was proposed, although the same
 3576 idea was covered by [ENOTTY], because the connotations of the name would
 3577 be misleading. It was pointed out that the *fcntl()* function uses the error code
 3578 [EINVAL] for this notion, and hence all instances of [EFTYPE] were changed
 3579 to this code.

3580 [EINTR] POSIX.1 prohibits conforming implementations from restarting interrupted
 3581 system calls of conforming applications unless the SA_RESTART flag is in
 3582 effect for the signal. However, it does not require that [EINTR] be returned
 3583 when another legitimate value may be substituted; for example, a partial
 3584 transfer count when *read()* or *write()* are interrupted. This is only given when
 3585 the signal-catching function returns normally as opposed to returns by
 3586 mechanisms like *longjmp()* or *siglongjmp()*.

3587 [ELOOP] In specifying conditions under which implementations would generate this
 3588 error, the following goals were considered:

- 3589 • To ensure that actual loops are detected, including loops that result from
 3590 symbolic links across distributed file systems.
- 3591 • To ensure that during pathname resolution an application can rely on the
 3592 ability to follow at least {SYMLOOP_MAX} symbolic links in the absence
 3593 of a loop.
- 3594 • To allow implementations to provide the capability of traversing more
 3595 than {SYMLOOP_MAX} symbolic links in the absence of a loop.
- 3596 • To allow implementations to detect loops and generate the error prior to
 3597 encountering {SYMLOOP_MAX} symbolic links.

3598	[ENAMETOOLONG]	
3599		When a symbolic link is encountered during pathname resolution, the
3600		contents of that symbolic link are used to create a new pathname. The
3601		standard developers intended to allow, but not require, that implementations
3602		enforce the restriction of {PATH_MAX} on the result of this pathname
3603		substitution.
3604	[ENOMEM]	The term “main memory” is not used in POSIX.1 because it is
3605		implementation-defined.
3606	[ENOTSUP]	This error code is to be used when an implementation chooses to implement
3607		the required functionality of IEEE Std 1003.1-2001 but does not support
3608		optional facilities defined by IEEE Std 1003.1-2001. The return of [ENOSYS] is
3609		to be taken to indicate that the function of the interface is not supported at all;
3610		the function will always fail with this error code.
3611	[ENOTTY]	The symbolic name for this error is derived from a time when device control
3612		was done by <i>ioctl()</i> and that operation was only permitted on a terminal
3613		interface. The term “TTY” is derived from “teletypewriter”, the devices to
3614		which this error originally applied.
3615	[EOVERFLOW]	Most of the uses of this error code are related to large file support. Typically,
3616		these cases occur on systems which support multiple programming
3617		environments with different sizes for <i>off_t</i> , but they may also occur in
3618		connection with remote file systems.
3619		In addition, when different programming environments have different widths
3620		for types such as <i>int</i> and <i>uid_t</i> , several functions may encounter a condition
3621		where a value in a particular environment is too wide to be represented. In
3622		that case, this error should be raised. For example, suppose the currently
3623		running process has 64-bit <i>int</i> , and file descriptor 9 223 372 036 854 775 807 is
3624		open and does not have the close-on-exec flag set. If the process then uses
3625		<i>execl()</i> to <i>exec</i> a file compiled in a programming environment with 32-bit <i>int</i> ,
3626		the call to <i>execl()</i> can fail with <i>errno</i> set to [EOVERFLOW]. A similar failure
3627		can occur with <i>execl()</i> if any of the user IDs or any of the group IDs to be
3628		assigned to the new process image are out of range for the executed file’s
3629		programming environment.
3630		Note, however, that this condition cannot occur for functions that are
3631		explicitly described as always being successful, such as <i>getpid()</i> .
3632	[EPIPE]	This condition normally generates the signal SIGPIPE; the error is returned if
3633		the signal does not terminate the process.
3634	[EROFS]	In historical implementations, attempting to <i>unlink()</i> or <i>rmdir()</i> a mount point
3635		would generate an [EBUSY] error. An implementation could be envisioned
3636		where such an operation could be performed without error. In this case, if
3637		<i>either</i> the directory entry or the actual data structures reside on a read-only file
3638		system, [EROFS] is the appropriate error to generate. (For example, changing
3639		the link count of a file on a read-only file system could not be done, as is
3640		required by <i>unlink()</i> , and thus an error should be reported.)
3641		Three error numbers, [EDOM], [EILSEQ], and [ERANGE], were added to this section primarily
3642		for consistency with the ISO C standard.

3643 **Alternative Solutions for Per-Thread *errno***

3644 The usual implementation of *errno* as a single global variable does not work in a multi-threaded
 3645 environment. In such an environment, a thread may make a POSIX.1 call and get a -1 error
 3646 return, but before that thread can check the value of *errno*, another thread might have made a
 3647 second POSIX.1 call that also set *errno*. This behavior is unacceptable in robust programs. There
 3648 were a number of alternatives that were considered for handling the *errno* problem:

- 3649 • Implement *errno* as a per-thread integer variable.
- 3650 • Implement *errno* as a service that can access the per-thread error number.
- 3651 • Change all POSIX.1 calls to accept an extra status argument and avoid setting *errno*.
- 3652 • Change all POSIX.1 calls to raise a language exception.

3653 The first option offers the highest level of compatibility with existing practice but requires
 3654 special support in the linker, compiler, and/or virtual memory system to support the new
 3655 concept of thread private variables. When compared with current practice, the third and fourth
 3656 options are much cleaner, more efficient, and encourage a more robust programming style, but
 3657 they require new versions of all of the POSIX.1 functions that might detect an error. The second
 3658 option offers compatibility with existing code that uses the `<errno.h>` header to define the
 3659 symbol *errno*. In this option, *errno* may be a macro defined:

```
3660     #define errno  (*__errno())
3661     extern int    *__errno();
```

3662 This option may be implemented as a per-thread variable whereby an *errno* field is allocated in
 3663 the user space object representing a thread, and whereby the function `__errno()` makes a system
 3664 call to determine the location of its user space object and returns the address of the *errno* field of
 3665 that object. Another implementation, one that avoids calling the kernel, involves allocating
 3666 stacks in chunks. The stack allocator keeps a side table indexed by chunk number containing a
 3667 pointer to the thread object that uses that chunk. The `__errno()` function then looks at the stack
 3668 pointer, determines the chunk number, and uses that as an index into the chunk table to find its
 3669 thread object and thus its private value of *errno*. On most architectures, this can be done in four
 3670 to five instructions. Some compilers may wish to implement `__errno()` inline to improve
 3671 performance.

3672 **Disallowing Return of the [EINTR] Error Code**

3673 Many blocking interfaces defined by IEEE Std 1003.1-2001 may return [EINTR] if interrupted
 3674 during their execution by a signal handler. Blocking interfaces introduced under the Threads
 3675 option do not have this property. Instead, they require that the interface appear to be atomic
 3676 with respect to interruption. In particular, clients of blocking interfaces need not handle any
 3677 possible [EINTR] return as a special case since it will never occur. If it is necessary to restart
 3678 operations or complete incomplete operations following the execution of a signal handler, this is
 3679 handled by the implementation, rather than by the application.

3680 Requiring applications to handle [EINTR] errors on blocking interfaces has been shown to be a
 3681 frequent source of often unreproducible bugs, and it adds no compelling value to the available
 3682 functionality. Thus, blocking interfaces introduced for use by multi-threaded programs do not
 3683 use this paradigm. In particular, in none of the functions `flockfile()`, `pthread_cond_timedwait()`,
 3684 `pthread_cond_wait()`, `pthread_join()`, `pthread_mutex_lock()`, and `sigwait()` did providing [EINTR]
 3685 returns add value, or even particularly make sense. Thus, these functions do not provide for an
 3686 [EINTR] return, even when interrupted by a signal handler. The same arguments can be applied
 3687 to `sem_wait()`, `sem_trywait()`, `sigwaitinfo()`, and `sigtimedwait()`, but implementations are
 3688 permitted to return [EINTR] error codes for these functions for compatibility with earlier

3689 versions of IEEE Std 1003.1. Applications cannot rely on calls to these functions returning
3690 [EINTR] error codes when signals are delivered to the calling thread, but they should allow for
3691 the possibility.

3692 B.2.3.1 Additional Error Numbers

3693 The ISO C standard defines the name space for implementations to add additional error
3694 numbers.

3695 B.2.4 Signal Concepts

3696 Historical implementations of signals, using the *signal()* function, have shortcomings that make
3697 them unreliable for many application uses. Because of this, a new signal mechanism, based very
3698 closely on the one of 4.2 BSD and 4.3 BSD, was added to POSIX.1.

3699 Signal Names

3700 The restriction on the actual type used for **sigset_t** is intended to guarantee that these objects can
3701 always be assigned, have their address taken, and be passed as parameters by value. It is not
3702 intended that this type be a structure including pointers to other data structures, as that could
3703 impact the portability of applications performing such operations. A reasonable implementation
3704 could be a structure containing an array of some integer type.

3705 The signals described in IEEE Std 1003.1-2001 must have unique values so that they may be
3706 named as parameters of **case** statements in the body of a C-language **switch** clause. However,
3707 implementation-defined signals may have values that overlap with each other or with signals
3708 specified in IEEE Std 1003.1-2001. An example of this is SIGABRT, which traditionally overlaps
3709 some other signal, such as SIGIOT.

3710 SIGKILL, SIGTERM, SIGUSR1, and SIGUSR2 are ordinarily generated only through the explicit
3711 use of the *kill()* function, although some implementations generate SIGKILL under
3712 extraordinary circumstances. SIGTERM is traditionally the default signal sent by the *kill*
3713 command.

3714 The signals SIGBUS, SIGEMT, SIGIOT, SIGTRAP, and SIGSYS were omitted from POSIX.1
3715 because their behavior is implementation-defined and could not be adequately categorized.
3716 Conforming implementations may deliver these signals, but must document the circumstances
3717 under which they are delivered and note any restrictions concerning their delivery. The signals
3718 SIGFPE, SIGILL, and SIGSEGV are similar in that they also generally result only from
3719 programming errors. They were included in POSIX.1 because they do indicate three relatively
3720 well-categorized conditions. They are all defined by the ISO C standard and thus would have to
3721 be defined by any system with an ISO C standard binding, even if not explicitly included in
3722 POSIX.1.

3723 There is very little that a Conforming POSIX.1 Application can do by catching, ignoring, or
3724 masking any of the signals SIGILL, SIGTRAP, SIGIOT, SIGEMT, SIGBUS, SIGSEGV, SIGSYS, or
3725 SIGFPE. They will generally be generated by the system only in cases of programming errors.
3726 While it may be desirable for some robust code (for example, a library routine) to be able to
3727 detect and recover from programming errors in other code, these signals are not nearly sufficient
3728 for that purpose. One portable use that does exist for these signals is that a command interpreter
3729 can recognize them as the cause of a process' termination (with *wait()*) and print an appropriate
3730 message. The mnemonic tags for these signals are derived from their PDP-11 origin.

3731 The signals SIGSTOP, SIGTSTP, SIGTTIN, SIGTTOU, and SIGCONT are provided for job control
3732 and are unchanged from 4.2 BSD. The signal SIGCHLD is also typically used by job control
3733 shells to detect children that have terminated or, as in 4.2 BSD, stopped.

3734 Some implementations, including System V, have a signal named SIGCLD, which is similar to
3735 SIGCHLD in 4.2 BSD. POSIX.1 permits implementations to have a single signal with both
3736 names. POSIX.1 carefully specifies ways in which conforming applications can avoid the
3737 semantic differences between the two different implementations. The name SIGCHLD was
3738 chosen for POSIX.1 because most current application usages of it can remain unchanged in
3739 conforming applications. SIGCLD in System V has more cases of semantics that POSIX.1 does
3740 not specify, and thus applications using it are more likely to require changes in addition to the
3741 name change.

3742 The signals SIGUSR1 and SIGUSR2 are commonly used by applications for notification of
3743 exceptional behavior and are described as “reserved as application-defined” so that such use is
3744 not prohibited. Implementations should not generate SIGUSR1 or SIGUSR2, except when
3745 explicitly requested by *kill()*. It is recommended that libraries not use these two signals, as such
3746 use in libraries could interfere with their use by applications calling the libraries. If such use is
3747 unavoidable, it should be documented. It is prudent for non-portable libraries to use non-
3748 standard signals to avoid conflicts with use of standard signals by portable libraries.

3749 There is no portable way for an application to catch or ignore non-standard signals. Some
3750 implementations define the range of signal numbers, so applications can install signal-catching
3751 functions for all of them. Unfortunately, implementation-defined signals often cause problems
3752 when caught or ignored by applications that do not understand the reason for the signal. While
3753 the desire exists for an application to be more robust by handling all possible signals (even those
3754 only generated by *kill()*), no existing mechanism was found to be sufficiently portable to include
3755 in POSIX.1. The value of such a mechanism, if included, would be diminished given that
3756 SIGKILL would still not be catchable.

3757 A number of new signal numbers are reserved for applications because the two user signals
3758 defined by POSIX.1 are insufficient for many realtime applications. A range of signal numbers is
3759 specified, rather than an enumeration of additional reserved signal names, because different
3760 applications and application profiles will require a different number of application signals. It is
3761 not desirable to burden all application domains and therefore all implementations with the
3762 maximum number of signals required by all possible applications. Note that in this context,
3763 signal numbers are essentially different signal priorities.

3764 The relatively small number of required additional signals, `{_POSIX_RTSIG_MAX}`, was chosen
3765 so as not to require an unreasonably large signal mask/set. While this number of signals defined
3766 in POSIX.1 will fit in a single 32-bit word signal mask, it is recognized that most existing
3767 implementations define many more signals than are specified in POSIX.1 and, in fact, many
3768 implementations have already exceeded 32 signals (including the “null signal”). Support of
3769 `{_POSIX_RTSIG_MAX}` additional signals may push some implementation over the single 32-bit
3770 word line, but is unlikely to push any implementations that are already over that line beyond the
3771 64-signal line.

3772 B.2.4.1 Signal Generation and Delivery

3773 The terms defined in this section are not used consistently in documentation of historical
3774 systems. Each signal can be considered to have a lifetime beginning with generation and ending
3775 with delivery or acceptance. The POSIX.1 definition of “delivery” does not exclude ignored
3776 signals; this is considered a more consistent definition. This revised text in several parts of
3777 IEEE Std 1003.1-2001 clarifies the distinct semantics of asynchronous signal delivery and
3778 synchronous signal acceptance. The previous wording attempted to categorize both under the
3779 term “delivery”, which led to conflicts over whether the effects of asynchronous signal delivery
3780 applied to synchronous signal acceptance.

3781 Signals generated for a process are delivered to only one thread. Thus, if more than one thread is
3782 eligible to receive a signal, one has to be chosen. The choice of threads is left entirely up to the
3783 implementation both to allow the widest possible range of conforming implementations and to
3784 give implementations the freedom to deliver the signal to the “easiest possible” thread should
3785 there be differences in ease of delivery between different threads.

3786 Note that should multiple delivery among cooperating threads be required by an application,
3787 this can be trivially constructed out of the provided single-delivery semantics. The construction
3788 of a *sigwait_multiple()* function that accomplishes this goal is presented with the rationale for
3789 *sigwaitinfo()*.

3790 Implementations should deliver unblocked signals as soon after they are generated as possible.
3791 However, it is difficult for POSIX.1 to make specific requirements about this, beyond those in
3792 *kill()* and *sigprocmask()*. Even on systems with prompt delivery, scheduling of higher priority
3793 processes is always likely to cause delays.

3794 In general, the interval between the generation and delivery of unblocked signals cannot be
3795 detected by an application. Thus, references to pending signals generally apply to blocked,
3796 pending signals. An implementation registers a signal as pending on the process when no thread
3797 has the signal unblocked and there are no threads blocked in a *sigwait()* function for that signal.
3798 Thereafter, the implementation delivers the signal to the first thread that unblocks the signal or
3799 calls a *sigwait()* function on a signal set containing this signal rather than choosing the recipient
3800 thread at the time the signal is sent.

3801 In the 4.3 BSD system, signals that are blocked and set to SIG_IGN are discarded immediately
3802 upon generation. For a signal that is ignored as its default action, if the action is SIG_DFL and
3803 the signal is blocked, a generated signal remains pending. In the 4.1 BSD system and in System V
3804 Release 3 (two other implementations that support a somewhat similar signal mechanism), all
3805 ignored blocked signals remain pending if generated. Because it is not normally useful for an
3806 application to simultaneously ignore and block the same signal, it was unnecessary for POSIX.1
3807 to specify behavior that would invalidate any of the historical implementations.

3808 There is one case in some historical implementations where an unblocked, pending signal does
3809 not remain pending until it is delivered. In the System V implementation of *signal()*, pending
3810 signals are discarded when the action is set to SIG_DFL or a signal-catching routine (as well as to
3811 SIG_IGN). Except in the case of setting SIGCHLD to SIG_DFL, implementations that do this do
3812 not conform completely to POSIX.1. Some earlier proposals for POSIX.1 explicitly stated this,
3813 but these statements were redundant due to the requirement that functions defined by POSIX.1
3814 not change attributes of processes defined by POSIX.1 except as explicitly stated.

3815 POSIX.1 specifically states that the order in which multiple, simultaneously pending signals are
3816 delivered is unspecified. This order has not been explicitly specified in historical
3817 implementations, but has remained quite consistent and been known to those familiar with the
3818 implementations. Thus, there have been cases where applications (usually system utilities) have
3819 been written with explicit or implicit dependencies on this order. Implementors and others
3820 porting existing applications may need to be aware of such dependencies.

3821 When there are multiple pending signals that are not blocked, implementations should arrange
3822 for the delivery of all signals at once, if possible. Some implementations stack calls to all pending
3823 signal-catching routines, making it appear that each signal-catcher was interrupted by the next
3824 signal. In this case, the implementation should ensure that this stacking of signals does not
3825 violate the semantics of the signal masks established by *sigaction()*. Other implementations
3826 process at most one signal when the operating system is entered, with remaining signals saved
3827 for later delivery. Although this practice is widespread, this behavior is neither standardized
3828 nor endorsed. In either case, implementations should attempt to deliver signals associated with
3829 the current state of the process (for example, SIGFPE) before other signals, if possible.

3830 In 4.2 BSD and 4.3 BSD, it is not permissible to ignore or explicitly block SIGCONT, because if
 3831 blocking or ignoring this signal prevented it from continuing a stopped process, such a process
 3832 could never be continued (only killed by SIGKILL). However, 4.2 BSD and 4.3 BSD do block
 3833 SIGCONT during execution of its signal-catching function when it is caught, creating exactly
 3834 this problem. A proposal was considered to disallow catching SIGCONT in addition to ignoring
 3835 and blocking it, but this limitation led to objections. The consensus was to require that
 3836 SIGCONT always continue a stopped process when generated. This removed the need to
 3837 disallow ignoring or explicit blocking of the signal; note that SIG_IGN and SIG_DFL are
 3838 equivalent for SIGCONT.

3839 B.2.4.2 Realtime Signal Generation and Delivery

3840 The Realtime Signals Extension option to POSIX.1 signal generation and delivery behavior is
 3841 required for the following reasons:

- 3842 • The **sigevent** structure is used by other POSIX.1 functions that result in asynchronous event
 3843 notifications to specify the notification mechanism to use and other information needed by
 3844 the notification mechanism. IEEE Std 1003.1-2001 defines only three symbolic values for the
 3845 notification mechanism. SIGEV_NONE is used to indicate that no notification is required
 3846 when the event occurs. This is useful for applications that use asynchronous I/O with polling
 3847 for completion. SIGEV_SIGNAL indicates that a signal is generated when the event occurs.
 3848 SIGEV_NOTIFY provides for “callback functions” for asynchronous notifications done by a
 3849 function call within the context of a new thread. This provides a multi-threaded process a
 3850 more natural means of notification than signals. The primary difficulty with previous
 3851 notification approaches has been to specify the environment of the notification routine.
 - 3852 — One approach is to limit the notification routine to call only functions permitted in a
 3853 signal handler. While the list of permissible functions is clearly stated, this is overly
 3854 restrictive.
 - 3855 — A second approach is to define a new list of functions or classes of functions that are
 3856 explicitly permitted or not permitted. This would give a programmer more lists to deal
 3857 with, which would be awkward.
 - 3858 — The third approach is to define completely the environment for execution of the
 3859 notification function. A clear definition of an execution environment for notification is
 3860 provided by executing the notification function in the environment of a newly created
 3861 thread.

3862 Implementations may support additional notification mechanisms by defining new values
 3863 for *sigev_notify*.

3864 For a notification type of SIGEV_SIGNAL, the other members of the **sigevent** structure
 3865 defined by IEEE Std 1003.1-2001 specify the realtime signal—that is, the signal number and
 3866 application-defined value that differentiates between occurrences of signals with the same
 3867 number—that will be generated when the event occurs. The structure is defined in
 3868 <signal.h>, even though the structure is not directly used by any of the signal functions,
 3869 because it is part of the signals interface used by the POSIX.1b “client functions”. When the
 3870 client functions include <signal.h> to define the signal names, the **sigevent** structure will
 3871 also be defined.

3872 An application-defined value passed to the signal handler is used to differentiate between
 3873 different “events” instead of requiring that the application use different signal numbers for
 3874 several reasons:

- 3875 — Realtime applications potentially handle a very large number of different events.
 3876 Requiring that implementations support a correspondingly large number of distinct

- 3877 signal numbers will adversely impact the performance of signal delivery because the
3878 signal masks to be manipulated on entry and exit to the handlers will become large.
- 3879 — Event notifications are prioritized by signal number (the rationale for this is explained in
3880 the following paragraphs) and the use of different signal numbers to differentiate
3881 between the different event notifications overloads the signal number more than has
3882 already been done. It also requires that the application writer make arbitrary assignments
3883 of priority to events that are logically of equal priority.
- 3884 A union is defined for the application-defined value so that either an integer constant or a
3885 pointer can be portably passed to the signal-catching function. On some architectures a
3886 pointer cannot be cast to an `int` and *vice versa*.
- 3887 Use of a structure here with an explicit notification type discriminant rather than explicit
3888 parameters to realtime functions, or embedded in other realtime structures, provides for
3889 future extensions to IEEE Std 1003.1-2001. Additional, perhaps more efficient, notification
3890 mechanisms can be supported for existing realtime function interfaces, such as timers and
3891 asynchronous I/O, by extending the `sigevent` structure appropriately. The existing realtime
3892 function interfaces will not have to be modified to use any such new notification mechanism.
3893 The revised text concerning the `SIGEV_SIGNAL` value makes consistent the semantics of the
3894 members of the `sigevent` structure, particularly in the definitions of `lio_listio()` and
3895 `aio_fsync()`. For uniformity, other revisions cause this specification to be referred to rather
3896 than inaccurately duplicated in the descriptions of functions and structures using the
3897 `sigevent` structure. The revised wording does not relax the requirement that the signal
3898 number be in the range `SIGRTMIN` to `SIGRTMAX` to guarantee queuing and passing of the
3899 application value, since that requirement is still implied by the signal names.
- 3900 • IEEE Std 1003.1-2001 is intentionally vague on whether “non-realtime” signal-generating
3901 mechanisms can result in a `siginfo_t` being supplied to the handler on delivery. In one
3902 existing implementation, a `siginfo_t` is posted on signal generation, even though the
3903 implementation does not support queuing of multiple occurrences of a signal. It is not the
3904 intent of IEEE Std 1003.1-2001 to preclude this, independent of the mandate to define signals
3905 that do support queuing. Any interpretation that appears to preclude this is a mistake in the
3906 reading or writing of the standard.
 - 3907 • Signals handled by realtime signal handlers might be generated by functions or conditions
3908 that do not allow the specification of an application-defined value and do not queue.
3909 IEEE Std 1003.1-2001 specifies the `si_code` member of the `siginfo_t` structure used in existing
3910 practice and defines additional codes so that applications can detect whether an application-
3911 defined value is present or not. The code `SI_USER` for `kill()`-generated signals is adopted
3912 from existing practice.
 - 3913 • The `sigaction()` `sa_flags` value `SA_SIGINFO` tells the implementation that the signal-catching
3914 function expects two additional arguments. When the flag is not set, a single argument, the
3915 signal number, is passed as specified by IEEE Std 1003.1-2001. Although IEEE Std 1003.1-2001
3916 does not explicitly allow the `info` argument to the handler function to be `NULL`, this is
3917 existing practice. This provides for compatibility with programs whose signal-catching
3918 functions are not prepared to accept the additional arguments. IEEE Std 1003.1-2001 is
3919 explicitly unspecified as to whether signals actually queue when `SA_SIGINFO` is not set for a
3920 signal, as there appear to be no benefits to applications in specifying one behavior or another.
3921 One existing implementation queues a `siginfo_t` on each signal generation, unless the signal
3922 is already pending, in which case the implementation discards the new `siginfo_t`; that is, the
3923 queue length is never greater than one. This implementation only examines `SA_SIGINFO` on
3924 signal delivery, discarding the queued `siginfo_t` if its delivery was not requested.

3925 IEEE Std 1003.1-2001 specifies several new values for the *si_code* member of the **siginfo_t**
 3926 structure. In existing practice, a *si_code* value of less than or equal to zero indicates that
 3927 the signal was generated by a process via the *kill()* function. In existing practice, values of *si_code*
 3928 that provide additional information for implementation-generated signals, such as SIGFPE or
 3929 SIGSEGV, are all positive. Thus, if implementations define the new constants specified in
 3930 IEEE Std 1003.1-2001 to be negative numbers, programs written to use existing practice will
 3931 not break. IEEE Std 1003.1-2001 chose not to attempt to specify existing practice values of
 3932 *si_code* other than SI_USER both because it was deemed beyond the scope of
 3933 IEEE Std 1003.1-2001 and because many of the values in existing practice appear to be
 3934 platform and implementation-defined. But, IEEE Std 1003.1-2001 does specify that if an
 3935 implementation—for example, one that does not have existing practice in this area—chooses
 3936 to define additional values for *si_code*, these values have to be different from the values of the
 3937 symbols specified by IEEE Std 1003.1-2001. This will allow conforming applications to
 3938 differentiate between signals generated by one of the POSIX.1b asynchronous events and
 3939 those generated by other implementation events in a manner compatible with existing
 3940 practice.

3941 The unique values of *si_code* for the POSIX.1b asynchronous events have implications for
 3942 implementations of, for example, asynchronous I/O or message passing in user space library
 3943 code. Such an implementation will be required to provide a hidden interface to the signal
 3944 generation mechanism that allows the library to specify the standard values of *si_code*.

3945 Existing practice also defines additional members of **siginfo_t**, such as the process ID and
 3946 user ID of the sending process for *kill()*-generated signals. These members were deemed not
 3947 necessary to meet the requirements of realtime applications and are not specified by
 3948 IEEE Std 1003.1-2001. Neither are they precluded.

3949 The third argument to the signal-catching function, *context*, is left undefined by
 3950 IEEE Std 1003.1-2001, but is specified in the interface because it matches existing practice for
 3951 the SA_SIGINFO flag. It was considered undesirable to require a separate implementation
 3952 for SA_SIGINFO for POSIX conformance on implementations that already support the two
 3953 additional parameters.

3954 • The requirement to deliver lower numbered signals in the range SIGRTMIN to SIGRTMAX
 3955 first, when multiple unblocked signals are pending, results from several considerations:

3956 — A method is required to prioritize event notifications. The signal number was chosen
 3957 instead of, for instance, associating a separate priority with each request, because an
 3958 implementation has to check pending signals at various points and select one for delivery
 3959 when more than one is pending. Specifying a selection order is the minimal additional
 3960 semantic that will achieve prioritized delivery. If a separate priority were to be associated
 3961 with queued signals, it would be necessary for an implementation to search all non-
 3962 empty, non-blocked signal queues and select from among them the pending signal with
 3963 the highest priority. This would significantly increase the cost of and decrease the
 3964 determinism of signal delivery.

3965 — Given the specified selection of the lowest numeric unblocked pending signal,
 3966 preemptive priority signal delivery can be achieved using signal numbers and signal
 3967 masks by ensuring that the *sa_mask* for each signal number blocks all signals with a
 3968 higher numeric value.

3969 For realtime applications that want to use only the newly defined realtime signal numbers
 3970 without interference from the standard signals, this can be achieved by blocking all of the
 3971 standard signals in the process signal mask and in the *sa_mask* installed by the signal
 3972 action for the realtime signal handlers.

3973 IEEE Std 1003.1-2001 explicitly leaves unspecified the ordering of signals outside of the range
 3974 of realtime signals and the ordering of signals within this range with respect to those outside
 3975 the range. It was believed that this would unduly constrain implementations or standards in
 3976 the future definition of new signals.

3977 B.2.4.3 Signal Actions

3978 Early proposals mentioned SIGCONT as a second exception to the rule that signals are not
 3979 delivered to stopped processes until continued. Because IEEE Std 1003.1-2001 now specifies that
 3980 SIGCONT causes the stopped process to continue when it is generated, delivery of SIGCONT is
 3981 not prevented because a process is stopped, even without an explicit exception to this rule.

3982 Ignoring a signal by setting the action to SIG_IGN (or SIG_DFL for signals whose default action
 3983 is to ignore) is not the same as installing a signal-catching function that simply returns. Invoking
 3984 such a function will interrupt certain system functions that block processes (for example, *wait()*,
 3985 *sigsuspend()*, *pause()*, *read()*, *write()*) while ignoring a signal has no such effect on the process.

3986 Historical implementations discard pending signals when the action is set to SIG_IGN.
 3987 However, they do not always do the same when the action is set to SIG_DFL and the default
 3988 action is to ignore the signal. IEEE Std 1003.1-2001 requires this for the sake of consistency and
 3989 also for completeness, since the only signal this applies to is SIGCHLD, and IEEE Std 1003.1-2001
 3990 disallows setting its action to SIG_IGN.

3991 Some implementations (System V, for example) assign different semantics for SIGCLD
 3992 depending on whether the action is set to SIG_IGN or SIG_DFL. Since POSIX.1 requires that the
 3993 default action for SIGCHLD be to ignore the signal, applications should always set the action to
 3994 SIG_DFL in order to avoid SIGCHLD.

3995 Whether or not an implementation allows SIG_IGN as a SIGCHLD disposition to be inherited
 3996 across a call to one of the *exec* family of functions or *posix_spawn()* is explicitly left as
 3997 unspecified. This change was made as a result of IEEE PASC Interpretation 1003.1 #132, and
 3998 permits the implementation to decide between the following alternatives:

- 3999 • Unconditionally leave SIGCHLD set to SIG_IGN, in which case the implementation would
 4000 not allow applications that assume inheritance of SIG_DFL to conform to
 4001 IEEE Std 1003.1-2001 without change. The implementation would, however, retain an ability
 4002 to control applications that create child processes but never call on the *wait* family of
 4003 functions, potentially filling up the process table.
- 4004 • Unconditionally reset SIGCHLD to SIG_DFL, in which case the implementation would allow
 4005 applications that assume inheritance of SIG_DFL to conform. The implementation would,
 4006 however, lose an ability to control applications that spawn child processes but never reap
 4007 them.
- 4008 • Provide some mechanism, not specified in IEEE Std 1003.1-2001, to control inherited
 4009 SIGCHLD dispositions.

4010 Some implementations (System V, for example) will deliver a SIGCLD signal immediately when
 4011 a process establishes a signal-catching function for SIGCLD when that process has a child that
 4012 has already terminated. Other implementations, such as 4.3 BSD, do not generate a new
 4013 SIGCHLD signal in this way. In general, a process should not attempt to alter the signal action
 4014 for the SIGCHLD signal while it has any outstanding children. However, it is not always
 4015 possible for a process to avoid this; for example, shells sometimes start up processes in pipelines
 4016 with other processes from the pipeline as children. Processes that cannot ensure that they have
 4017 no children when altering the signal action for SIGCHLD thus need to be prepared for, but not
 4018 depend on, generation of an immediate SIGCHLD signal.

4019 The default action of the stop signals (SIGSTOP, SIGTSTP, SIGTTIN, SIGTTOU) is to stop a
 4020 process that is executing. If a stop signal is delivered to a process that is already stopped, it has
 4021 no effect. In fact, if a stop signal is generated for a stopped process whose signal mask blocks the
 4022 signal, the signal will never be delivered to the process since the process must receive a
 4023 SIGCONT, which discards all pending stop signals, in order to continue executing.

4024 The SIGCONT signal continues a stopped process even if SIGCONT is blocked (or ignored).
 4025 However, if a signal-catching routine has been established for SIGCONT, it will not be entered
 4026 until SIGCONT is unblocked.

4027 If a process in an orphaned process group stops, it is no longer under the control of a job control
 4028 shell and hence would not normally ever be continued. Because of this, orphaned processes that
 4029 receive terminal-related stop signals (SIGTSTP, SIGTTIN, SIGTTOU, but not SIGSTOP) must not
 4030 be allowed to stop. The goal is to prevent stopped processes from languishing forever. (As
 4031 SIGSTOP is sent only via *kill()*, it is assumed that the process or user sending a SIGSTOP can
 4032 send a SIGCONT when desired.) Instead, the system must discard the stop signal. As an
 4033 extension, it may also deliver another signal in its place. 4.3 BSD sends a SIGKILL, which is
 4034 overly effective because SIGKILL is not catchable. Another possible choice is SIGHUP. 4.3 BSD
 4035 also does this for orphaned processes (processes whose parent has terminated) rather than for
 4036 members of orphaned process groups; this is less desirable because job control shells manage
 4037 process groups. POSIX.1 also prevents SIGTTIN and SIGTTOU signals from being generated for
 4038 processes in orphaned process groups as a direct result of activity on a terminal, preventing
 4039 infinite loops when *read()* and *write()* calls generate signals that are discarded; see Section
 4040 A.11.1.4 (on page 68). A similar restriction on the generation of SIGTSTP was considered, but
 4041 that would be unnecessary and more difficult to implement due to its asynchronous nature.

4042 Although POSIX.1 requires that signal-catching functions be called with only one argument,
 4043 there is nothing to prevent conforming implementations from extending POSIX.1 to pass
 4044 additional arguments, as long as Strictly Conforming POSIX.1 Applications continue to compile
 4045 and execute correctly. Most historical implementations do, in fact, pass additional, signal-
 4046 specific arguments to certain signal-catching routines.

4047 There was a proposal to change the declared type of the signal handler to:

```
4048     void func (int sig, ...);
```

4049 The usage of ellipses ("...") is ISO C standard syntax to indicate a variable number of
 4050 arguments. Its use was intended to allow the implementation to pass additional information to
 4051 the signal handler in a standard manner.

4052 Unfortunately, this construct would require all signal handlers to be defined with this syntax
 4053 because the ISO C standard allows implementations to use a different parameter passing
 4054 mechanism for variable parameter lists than for non-variable parameter lists. Thus, all existing
 4055 signal handlers in all existing applications would have to be changed to use the variable syntax
 4056 in order to be standard and portable. This is in conflict with the goal of Minimal Changes to
 4057 Existing Application Code.

4058 When terminating a process from a signal-catching function, processes should be aware of any
 4059 interpretation that their parent may make of the status returned by *wait()* or *waitpid()*. In
 4060 particular, a signal-catching function should not call *exit(0)* or *_exit(0)* unless it wants to indicate
 4061 successful termination. A non-zero argument to *exit()* or *_exit()* can be used to indicate
 4062 unsuccessful termination. Alternatively, the process can use *kill()* to send itself a fatal signal
 4063 (first ensuring that the signal is set to the default action and not blocked). See also the
 4064 RATIONALE section of the *_exit()* function.

4065 The behavior of *unsafe* functions, as defined by this section, is undefined when they are invoked
 4066 from signal-catching functions in certain circumstances. The behavior of reentrant functions, as

4067 defined by this section, is as specified by POSIX.1, regardless of invocation from a signal-
4068 catching function. This is the only intended meaning of the statement that reentrant functions
4069 may be used in signal-catching functions without restriction. Applications must still consider all
4070 effects of such functions on such things as data structures, files, and process state. In particular,
4071 application writers need to consider the restrictions on interactions when interrupting *sleep()*
4072 (see *sleep()*) and interactions among multiple handles for a file description. The fact that any
4073 specific function is listed as reentrant does not necessarily mean that invocation of that function
4074 from a signal-catching function is recommended.

4075 In order to prevent errors arising from interrupting non-reentrant function calls, applications
4076 should protect calls to these functions either by blocking the appropriate signals or through the
4077 use of some programmatic semaphore. POSIX.1 does not address the more general problem of
4078 synchronizing access to shared data structures. Note in particular that even the “safe” functions
4079 may modify the global variable *errno*; the signal-catching function may want to save and restore
4080 its value. The same principles apply to the reentrancy of application routines and asynchronous
4081 data access.

4082 Note that *longjmp()* and *siglongjmp()* are not in the list of reentrant functions. This is because the
4083 code executing after *longjmp()* or *siglongjmp()* can call any unsafe functions with the same
4084 danger as calling those unsafe functions directly from the signal handler. Applications that use
4085 *longjmp()* or *siglongjmp()* out of signal handlers require rigorous protection in order to be
4086 portable. Many of the other functions that are excluded from the list are traditionally
4087 implemented using either the C language *malloc()* or *free()* functions or the ISO C standard I/O
4088 library, both of which traditionally use data structures in a non-reentrant manner. Because any
4089 combination of different functions using a common data structure can cause reentrancy
4090 problems, POSIX.1 does not define the behavior when any unsafe function is called in a signal
4091 handler that interrupts any unsafe function.

4092 The only realtime extension to signal actions is the addition of the additional parameters to the
4093 signal-catching function. This extension has been explained and motivated in the previous
4094 section. In making this extension, though, developers of POSIX.1b ran into issues relating to
4095 function prototypes. In response to input from the POSIX.1 standard developers, members were
4096 added to the **sigaction** structure to specify function prototypes for the newer signal-catching
4097 function specified by POSIX.1b. These members follow changes that are being made to POSIX.1.
4098 Note that IEEE Std 1003.1-2001 explicitly states that these fields may overlap so that a union can
4099 be defined. This enabled existing implementations of POSIX.1 to maintain binary-compatibility
4100 when these extensions were added.

4101 The **siginfo_t** structure was adopted for passing the application-defined value to match existing
4102 practice, but the existing practice has no provision for an application-defined value, so this was
4103 added. Note that POSIX normally reserves the “_t” type designation for opaque types. The
4104 **siginfo_t** structure breaks with this convention to follow existing practice and thus promote
4105 portability. Standardization of the existing practice for the other members of this structure may
4106 be addressed in the future.

4107 Although it is not explicitly visible to applications, there are additional semantics for signal
4108 actions implied by queued signals and their interaction with other POSIX.1b realtime functions.
4109 Specifically:

- 4110 • It is not necessary to queue signals whose action is SIG_IGN.
- 4111 • For implementations that support POSIX.1b timers, some interaction with the timer functions
4112 at signal delivery is implied to manage the timer overrun count.

4113 **B.2.4.4** *Signal Effects on Other Functions*

4114 The most common behavior of an interrupted function after a signal-catching function returns is
 4115 for the interrupted function to give an [EINTR] error unless the SA_RESTART flag is in effect for
 4116 the signal. However, there are a number of specific exceptions, including *sleep()* and certain
 4117 situations with *read()* and *write()*.

4118 The historical implementations of many functions defined by IEEE Std 1003.1-2001 are not
 4119 interruptible, but delay delivery of signals generated during their execution until after they
 4120 complete. This is never a problem for functions that are guaranteed to complete in a short
 4121 (imperceptible to a human) period of time. It is normally those functions that can suspend a
 4122 process indefinitely or for long periods of time (for example, *wait()*, *pause()*, *sigsuspend()*, *sleep()*,
 4123 or *read()/write()* on a slow device like a terminal) that are interruptible. This permits
 4124 applications to respond to interactive signals or to set timeouts on calls to most such functions
 4125 with *alarm()*. Therefore, implementations should generally make such functions (including ones
 4126 defined as extensions) interruptible.

4127 Functions not mentioned explicitly as interruptible may be so on some implementations,
 4128 possibly as an extension where the function gives an [EINTR] error. There are several functions
 4129 (for example, *getpid()*, *getuid()*) that are specified as never returning an error, which can thus
 4130 never be extended in this way.

4131 If a signal-catching function returns while the SA_RESTART flag is in effect, an interrupted
 4132 function is restarted at the point it was interrupted. Conforming applications cannot make
 4133 assumptions about the internal behavior of interrupted functions, even if the functions are
 4134 async-signal-safe. For example, suppose the *read()* function is interrupted with SA_RESTART in
 4135 effect, the signal-catching function closes the file descriptor being read from and returns, and the
 4136 *read()* function is then restarted; in this case the application cannot assume that the *read()*
 4137 function will give an [EBADF] error, since *read()* might have checked the file descriptor for
 4138 validity before being interrupted.

4139 **B.2.5** **Standard I/O Streams**4140 **B.2.5.1** *Interaction of File Descriptors and Standard I/O Streams*

4141 There is no additional rationale provided for this section.

4142 **B.2.5.2** *Stream Orientation and Encoding Rules*

4143 There is no additional rationale provided for this section.

4144 **B.2.6** **STREAMS**

4145 STREAMS are introduced into IEEE Std 1003.1-2001 as part of the alignment with the Single
 4146 UNIX Specification, but marked as an option in recognition that not all systems may wish to
 4147 implement the facility. The option within IEEE Std 1003.1-2001 is denoted by the XSR margin
 4148 marker. The standard developers made this option independent of the XSI option.

4149 STREAMS are a method of implementing network services and other character-based
 4150 input/output mechanisms, with the STREAM being a full-duplex connection between a process
 4151 and a device. STREAMS provides direct access to protocol modules, and optional protocol
 4152 modules can be interposed between the process-end of the STREAM and the device-driver at the
 4153 device-end of the STREAM. Pipes can be implemented using the STREAMS mechanism, so they
 4154 can provide process-to-process as well as process-to-device communications.

4155 This section introduces STREAMS I/O, the message types used to control them, an overview of
4156 the priority mechanism, and the interfaces used to access them.

4157 *B.2.6.1 Accessing STREAMS*

4158 There is no additional rationale provided for this section.

4159 **B.2.7 XSI Interprocess Communication**

4160 There are two forms of IPC supported as options in IEEE Std 1003.1-2001. The traditional
4161 System V IPC routines derived from the SVID—that is, the *msg**(), *sem**(), and *shm**()
4162 interfaces—are mandatory on XSI-conformant systems. Thus, all XSI-conformant systems
4163 provide the same mechanisms for manipulating messages, shared memory, and semaphores.

4164 In addition, the POSIX Realtime Extension provides an alternate set of routines for those systems
4165 supporting the appropriate options.

4166 The application writer is presented with a choice: the System V interfaces or the POSIX
4167 interfaces (loosely derived from the Berkeley interfaces). The XSI profile prefers the System V
4168 interfaces, but the POSIX interfaces may be more suitable for realtime or other performance-
4169 sensitive applications.

4170 *B.2.7.1 IPC General Information*

4171 General information that is shared by all three mechanisms is described in this section. The
4172 common permissions mechanism is briefly introduced, describing the mode bits, and how they
4173 are used to determine whether or not a process has access to read or write/alter the appropriate
4174 instance of one of the IPC mechanisms. All other relevant information is contained in the
4175 reference pages themselves.

4176 The semaphore type of IPC allows processes to communicate through the exchange of
4177 semaphore values. A semaphore is a positive integer. Since many applications require the use of
4178 more than one semaphore, XSI-conformant systems have the ability to create sets or arrays of
4179 semaphores.

4180 Calls to support semaphores include:

4181 *semctl()*, *semget()*, *semop()*

4182 Semaphore sets are created by using the *semget()* function.

4183 The message type of IPC allows processes to communicate through the exchange of data stored
4184 in buffers. This data is transmitted between processes in discrete portions known as messages.

4185 Calls to support message queues include:

4186 *msgctl()*, *msgget()*, *msgrcv()*, *msgsnd()*

4187 The shared memory type of IPC allows two or more processes to share memory and
4188 consequently the data contained therein. This is done by allowing processes to set up access to a
4189 common memory address space. This sharing of memory provides a fast means of exchange of
4190 data between processes.

4191 Calls to support shared memory include:

4192 *shmctl()*, *shmdt()*, *shmget()*

4193 The *ftok()* interface is also provided.

4194 **B.2.8 Realtime**4195 **Advisory Information**

4196 POSIX.1b contains an Informative Annex with proposed interfaces for “realtime files”. These
 4197 interfaces could determine groups of the exact parameters required to do “direct I/O” or
 4198 “extents”. These interfaces were objected to by a significant portion of the balloting group as too
 4199 complex. A conforming application had little chance of correctly navigating the large parameter
 4200 space to match its desires to the system. In addition, they only applied to a new type of file
 4201 (realtime files) and they told the implementation exactly what to do as opposed to advising the
 4202 implementation on application behavior and letting it optimize for the system the (portable)
 4203 application was running on. For example, it was not clear how a system that had a disk array
 4204 should set its parameters.

4205 There seemed to be several overall goals:

- 4206 • Optimizing sequential access
- 4207 • Optimizing caching behavior
- 4208 • Optimizing I/O data transfer
- 4209 • Preallocation

4210 The advisory interfaces, *posix_fadvise()* and *posix_madvise()*, satisfy the first two goals. The
 4211 `POSIX_FADV_SEQUENTIAL` and `POSIX_MADV_SEQUENTIAL` advice tells the
 4212 implementation to expect serial access. Typically the system will prefetch the next several serial
 4213 accesses in order to overlap I/O. It may also free previously accessed serial data if memory is
 4214 tight. If the application is not doing serial access it can use `POSIX_FADV_WILLNEED` and
 4215 `POSIX_MADV_WILLNEED` to accomplish I/O overlap, as required. When the application
 4216 advises `POSIX_FADV_RANDOM` or `POSIX_MADV_RANDOM` behavior, the implementation
 4217 usually tries to fetch a minimum amount of data with each request and it does not expect much
 4218 locality. `POSIX_FADV_DONTNEED` and `POSIX_MADV_DONTNEED` allow the system to free
 4219 up caching resources as the data will not be required in the near future.

4220 `POSIX_FADV_NOREUSE` tells the system that caching the specified data is not optimal. For file
 4221 I/O, the transfer should go directly to the user buffer instead of being cached internally by the
 4222 implementation. To portably perform direct disk I/O on all systems, the application must
 4223 perform its I/O transfers according to the following rules:

- 4224 1. The user buffer should be aligned according to the `{POSIX_REC_XFER_ALIGN}` *pathconf()*
 4225 variable.
- 4226 2. The number of bytes transferred in an I/O operation should be a multiple of the
 4227 `{POSIX_ALLOC_SIZE_MIN}` *pathconf()* variable.
- 4228 3. The offset into the file at the start of an I/O operation should be a multiple of the
 4229 `{POSIX_ALLOC_SIZE_MIN}` *pathconf()* variable.
- 4230 4. The application should ensure that all threads which open a given file specify
 4231 `POSIX_FADV_NOREUSE` to be sure that there is no unexpected interaction between
 4232 threads using buffered I/O and threads using direct I/O to the same file.

4233 In some cases, a user buffer must be properly aligned in order to be transferred directly to/from
 4234 the device. The `{POSIX_REC_XFER_ALIGN}` *pathconf()* variable tells the application the proper
 4235 alignment.

4236 The preallocation goal is met by the space control function, *posix_fallocate()*. The application can
 4237 use *posix_fallocate()* to guarantee no `[ENOSPC]` errors and to improve performance by prepaying

4238 any overhead required for block allocation.

4239 Implementations may use information conveyed by a previous *posix_fadvise()* call to influence
4240 the manner in which allocation is performed. For example, if an application did the following
4241 calls:

```
4242     fd = open("file");
4243     posix_fadvise(fd, offset, len, POSIX_FADV_SEQUENTIAL);
4244     posix_fallocate(fd, len, size);
```

4245 an implementation might allocate the file contiguously on disk.

4246 Finally, the *pathconf()* variables {*POSIX_REC_MIN_XFER_SIZE*},
4247 {*POSIX_REC_MAX_XFER_SIZE*}, and {*POSIX_REC_INCR_XFER_SIZE*} tell the application a
4248 range of transfer sizes that are recommended for best I/O performance.

4249 Where bounded response time is required, the vendor can supply the appropriate settings of the
4250 advisories to achieve a guaranteed performance level.

4251 The interfaces meet the goals while allowing applications using regular files to take advantage of
4252 performance optimizations. The interfaces tell the implementation expected application
4253 behavior which the implementation can use to optimize performance on a particular system
4254 with a particular dynamic load.

4255 The *posix_memalign()* function was added to allow for the allocation of specifically aligned
4256 buffers; for example, for {*POSIX_REC_XFER_ALIGN*}.

4257 The working group also considered the alternative of adding a function which would return an
4258 aligned pointer to memory within a user-supplied buffer. This was not considered to be the best
4259 method, because it potentially wastes large amounts of memory when buffers need to be aligned
4260 on large alignment boundaries.

4261 **Message Passing**

4262 This section provides the rationale for the definition of the message passing interface in
4263 IEEE Std 1003.1-2001. This is presented in terms of the objectives, models, and requirements
4264 imposed upon this interface.

4265 • Objectives

4266 Many applications, including both realtime and database applications, require a means of
4267 passing arbitrary amounts of data between cooperating processes comprising the overall
4268 application on one or more processors. Many conventional interfaces for interprocess
4269 communication are insufficient for realtime applications in that efficient and deterministic
4270 data passing methods cannot be implemented. This has prompted the definition of message
4271 passing interfaces providing these facilities:

- 4272 — Open a message queue.
- 4273 — Send a message to a message queue.
- 4274 — Receive a message from a queue, either synchronously or asynchronously.
- 4275 — Alter message queue attributes for flow and resource control.

4276 It is assumed that an application may consist of multiple cooperating processes and that
4277 these processes may wish to communicate and coordinate their activities. The message
4278 passing facility described in IEEE Std 1003.1-2001 allows processes to communicate through
4279 system-wide queues. These message queues are accessed through names that may be
4280 pathnames. A message queue can be opened for use by multiple sending and/or multiple

- 4281 receiving processes.
- 4282 • Background on Embedded Applications
- 4283 Interprocess communication utilizing message passing is a key facility for the construction of
4284 deterministic, high-performance realtime applications. The facility is present in all realtime
4285 systems and is the framework upon which the application is constructed. The performance of
4286 the facility is usually a direct indication of the performance of the resulting application.
- 4287 Realtime applications, especially for embedded systems, are typically designed around the
4288 performance constraints imposed by the message passing mechanisms. Applications for
4289 embedded systems are typically very tightly constrained. Application writers expect to
4290 design and control the entire system. In order to minimize system costs, the writer will
4291 attempt to use all resources to their utmost and minimize the requirement to add additional
4292 memory or processors.
- 4293 The embedded applications usually share address spaces and only a simple message passing
4294 mechanism is required. The application can readily access common data incurring only
4295 mutual-exclusion overheads. The models desired are the simplest possible with the
4296 application building higher-level facilities only when needed.
- 4297 • Requirements
- 4298 The following requirements determined the features of the message passing facilities defined
4299 in IEEE Std 1003.1-2001:
- 4300 — Naming of Message Queues
- 4301 The mechanism for gaining access to a message queue is a pathname evaluated in a
4302 context that is allowed to be a file system name space, or it can be independent of any file
4303 system. This is a specific attempt to allow implementations based on either method in
4304 order to address both embedded systems and to also allow implementation in larger
4305 systems.
- 4306 The interface of *mq_open()* is defined to allow but not require the access control and name
4307 conflicts resulting from utilizing a file system for name resolution. All required behavior
4308 is specified for the access control case. Yet a conforming implementation, such as an
4309 embedded system kernel, may define that there are no distinctions between users and
4310 may define that all processes have all access privileges.
- 4311 — Embedded System Naming
- 4312 Embedded systems need to be able to utilize independent name spaces for accessing the
4313 various system objects. They typically do not have a file system, precluding its utilization
4314 as a common name resolution mechanism. The modularity of an embedded system limits
4315 the connections between separate mechanisms that can be allowed.
- 4316 Embedded systems typically do not have any access protection. Since the system does not
4317 support the mixing of applications from different areas, and usually does not even have
4318 the concept of an authorization entity, access control is not useful.
- 4319 — Large System Naming
- 4320 On systems with more functionality, the name resolution must support the ability to use
4321 the file system as the name resolution mechanism/object storage medium and to have
4322 control over access to the objects. Utilizing the pathname space can result in further errors
4323 when the names conflict with other objects.

4324 — Fixed Size of Messages

4325 The interfaces impose a fixed upper bound on the size of messages that can be sent to a
 4326 specific message queue. The size is set on an individual queue basis and cannot be
 4327 changed dynamically.

4328 The purpose of the fixed size is to increase the ability of the system to optimize the
 4329 implementation of *mq_send()* and *mq_receive()*. With fixed sizes of messages and fixed
 4330 numbers of messages, specific message blocks can be pre-allocated. This eliminates a
 4331 significant amount of checking for errors and boundary conditions. Additionally, an
 4332 implementation can optimize data copying to maximize performance. Finally, with a
 4333 restricted range of message sizes, an implementation is better able to provide
 4334 deterministic operations.

4335 — Prioritization of Messages

4336 Message prioritization allows the application to determine the order in which messages
 4337 are received. Prioritization of messages is a key facility that is provided by most realtime
 4338 kernels and is heavily utilized by the applications. The major purpose of having priorities
 4339 in message queues is to avoid priority inversions in the message system, where a high-
 4340 priority message is delayed behind one or more lower-priority messages. This allows the
 4341 applications to be designed so that they do not need to be interrupted in order to change
 4342 the flow of control when exceptional conditions occur. The prioritization does add
 4343 additional overhead to the message operations in those cases it is actually used but a
 4344 clever implementation can optimize for the FIFO case to make that more efficient.

4345 — Asynchronous Notification

4346 The interface supports the ability to have a task asynchronously notified of the
 4347 availability of a message on the queue. The purpose of this facility is to allow the task to
 4348 perform other functions and yet still be notified that a message has become available on
 4349 the queue.

4350 To understand the requirement for this function, it is useful to understand two models of
 4351 application design: a single task performing multiple functions and multiple tasks
 4352 performing a single function. Each of these models has advantages.

4353 Asynchronous notification is required to build the model of a single task performing
 4354 multiple operations. This model typically results from either the expectation that
 4355 interruption is less expensive than utilizing a separate task or from the growth of the
 4356 application to include additional functions.

4357 **Semaphores**

4358 Semaphores are a high-performance process synchronization mechanism. Semaphores are
 4359 named by null-terminated strings of characters.

4360 A semaphore is created using the *sem_init()* function or the *sem_open()* function with the
 4361 *O_CREAT* flag set in *oflag*.

4362 To use a semaphore, a process has to first initialize the semaphore or inherit an open descriptor
 4363 for the semaphore via *fork()*.

4364 A semaphore preserves its state when the last reference is closed. For example, if a semaphore
 4365 has a value of 13 when the last reference is closed, it will have a value of 13 when it is next
 4366 opened.

4367 When a semaphore is created, an initial state for the semaphore has to be provided. This value is
 4368 a non-negative integer. Negative values are not possible since they indicate the presence of

4369 blocked processes. The persistence of any of these objects across a system crash or a system
4370 reboot is undefined. Conforming applications must not depend on any sort of persistence across
4371 a system reboot or a system crash.

4372 • Models and Requirements

4373 A realtime system requires synchronization and communication between the processes
4374 comprising the overall application. An efficient and reliable synchronization mechanism has
4375 to be provided in a realtime system that will allow more than one schedulable process
4376 mutually-exclusive access to the same resource. This synchronization mechanism has to
4377 allow for the optimal implementation of synchronization or systems implementors will
4378 define other, more cost-effective methods.

4379 At issue are the methods whereby multiple processes (tasks) can be designed and
4380 implemented to work together in order to perform a single function. This requires
4381 interprocess communication and synchronization. A semaphore mechanism is the lowest
4382 level of synchronization that can be provided by an operating system.

4383 A semaphore is defined as an object that has an integral value and a set of blocked processes
4384 associated with it. If the value is positive or zero, then the set of blocked processes is empty;
4385 otherwise, the size of the set is equal to the absolute value of the semaphore value. The value
4386 of the semaphore can be incremented or decremented by any process with access to the
4387 semaphore and must be done as an indivisible operation. When a semaphore value is less
4388 than or equal to zero, any process that attempts to lock it again will block or be informed that
4389 it is not possible to perform the operation.

4390 A semaphore may be used to guard access to any resource accessible by more than one
4391 schedulable task in the system. It is a global entity and not associated with any particular
4392 process. As such, a method of obtaining access to the semaphore has to be provided by the
4393 operating system. A process that wants access to a critical resource (section) has to wait on
4394 the semaphore that guards that resource. When the semaphore is locked on behalf of a
4395 process, it knows that it can utilize the resource without interference by any other
4396 cooperating process in the system. When the process finishes its operation on the resource,
4397 leaving it in a well-defined state, it posts the semaphore, indicating that some other process
4398 may now obtain the resource associated with that semaphore.

4399 In this section, mutexes and condition variables are specified as the synchronization
4400 mechanisms between threads.

4401 These primitives are typically used for synchronizing threads that share memory in a single
4402 process. However, this section provides an option allowing the use of these synchronization
4403 interfaces and objects between processes that share memory, regardless of the method for
4404 sharing memory.

4405 Much experience with semaphores shows that there are two distinct uses of synchronization:
4406 locking, which is typically of short duration; and waiting, which is typically of long or
4407 unbounded duration. These distinct usages map directly onto mutexes and condition
4408 variables, respectively.

4409 Semaphores are provided in IEEE Std 1003.1-2001 primarily to provide a means of
4410 synchronization for processes; these processes may or may not share memory. Mutexes and
4411 condition variables are specified as synchronization mechanisms between threads; these
4412 threads always share (some) memory. Both are synchronization paradigms that have been in
4413 widespread use for a number of years. Each set of primitives is particularly well matched to
4414 certain problems.

4415 With respect to binary semaphores, experience has shown that condition variables and
4416 mutexes are easier to use for many synchronization problems than binary semaphores. The
4417 primary reason for this is the explicit appearance of a Boolean predicate that specifies when
4418 the condition wait is satisfied. This Boolean predicate terminates a loop, including the call to
4419 *pthread_cond_wait()*. As a result, extra wakeups are benign since the predicate governs
4420 whether the thread will actually proceed past the condition wait. With stateful primitives,
4421 such as binary semaphores, the wakeup in itself typically means that the wait is satisfied. The
4422 burden of ensuring correctness for such waits is thus placed on *all* signalers of the semaphore
4423 rather than on an *explicitly coded* Boolean predicate located at the condition wait. Experience
4424 has shown that the latter creates a major improvement in safety and ease-of-use.

4425 Counting semaphores are well matched to dealing with producer/consumer problems,
4426 including those that might exist between threads of different processes, or between a signal
4427 handler and a thread. In the former case, there may be little or no memory shared by the
4428 processes; in the latter case, one is not communicating between co-equal threads, but
4429 between a thread and an interrupt-like entity. It is for these reasons that IEEE Std 1003.1-2001
4430 allows semaphores to be used by threads.

4431 Mutexes and condition variables have been effectively used with and without priority
4432 inheritance, priority ceiling, and other attributes to synchronize threads that share memory.
4433 The efficiency of their implementation is comparable to or better than that of other
4434 synchronization primitives that are sometimes harder to use (for example, binary
4435 semaphores). Furthermore, there is at least one known implementation of Ada tasking that
4436 uses these primitives. Mutexes and condition variables together constitute an appropriate,
4437 sufficient, and complete set of inter-thread synchronization primitives.

4438 Efficient multi-threaded applications require high-performance synchronization primitives.
4439 Considerations of efficiency and generality require a small set of primitives upon which more
4440 sophisticated synchronization functions can be built.

4441 • Standardization Issues

4442 It is possible to implement very high-performance semaphores using test-and-set
4443 instructions on shared memory locations. The library routines that implement such a high-
4444 performance interface have to properly ensure that a *sem_wait()* or *sem_trywait()*
4445 operation that cannot be performed will issue a blocking semaphore system call or properly report the
4446 condition to the application. The same interface to the application program would be
4447 provided by a high-performance implementation.

4448 *B.2.8.1 Realtime Signals*

4449 **Realtime Signals Extension**

4450 This portion of the rationale presents models, requirements, and standardization issues relevant
4451 to the Realtime Signals Extension. This extension provides the capability required to support
4452 reliable, deterministic, asynchronous notification of events. While a new mechanism,
4453 unencumbered by the historical usage and semantics of POSIX.1 signals, might allow for a more
4454 efficient implementation, the application requirements for event notification can be met with a
4455 small number of extensions to signals. Therefore, a minimal set of extensions to signals to
4456 support the application requirements is specified.

4457 The realtime signal extensions specified in this section are used by other realtime functions
4458 requiring asynchronous notification:

4459 • Models

4460 The model supported is one of multiple cooperating processes, each of which handles
 4461 multiple asynchronous external events. Events represent occurrences that are generated as
 4462 the result of some activity in the system. Examples of occurrences that can constitute an
 4463 event include:

- 4464 — Completion of an asynchronous I/O request
- 4465 — Expiration of a POSIX.1b timer
- 4466 — Arrival of an interprocess message
- 4467 — Generation of a user-defined event

4468 Processing of these events may occur synchronously via polling for event notifications or
 4469 asynchronously via a software interrupt mechanism. Existing practice for this model is well
 4470 established for traditional proprietary realtime operating systems, realtime executives, and
 4471 realtime extended POSIX-like systems.

4472 A contrasting model is that of “cooperating sequential processes” where each process
 4473 handles a single priority of events via polling. Each process blocks while waiting for events,
 4474 and each process depends on the preemptive, priority-based process scheduling mechanism
 4475 to arbitrate between events of different priority that need to be processed concurrently.
 4476 Existing practice for this model is also well established for small realtime executives that
 4477 typically execute in an unprotected physical address space, but it is just emerging in the
 4478 context of a fuller function operating system with multiple virtual address spaces.

4479 It could be argued that the cooperating sequential process model, and the facilities supported
 4480 by the POSIX Threads Extension obviate a software interrupt model. But, even with the
 4481 cooperating sequential process model, the need has been recognized for a software interrupt
 4482 model to handle exceptional conditions and process aborting, so the mechanism must be
 4483 supported in any case. Furthermore, it is not the purview of IEEE Std 1003.1-2001 to attempt
 4484 to convince realtime practitioners that their current application models based on software
 4485 interrupts are “broken” and should be replaced by the cooperating sequential process model.
 4486 Rather, it is the charter of IEEE Std 1003.1-2001 to provide standard extensions to
 4487 mechanisms that support existing realtime practice.

4488 • Requirements

4489 This section discusses the following realtime application requirements for asynchronous
 4490 event notification:

- 4491 — Reliable delivery of asynchronous event notification

4492 The events notification mechanism guarantees delivery of an event notification.
 4493 Asynchronous operations (such as asynchronous I/O and timers) that complete
 4494 significantly after they are invoked have to guarantee that delivery of the event
 4495 notification can occur at the time of completion.

- 4496 — Prioritized handling of asynchronous event notifications

4497 The events notification mechanism supports the assigning of a user function as an event
 4498 notification handler. Furthermore, the mechanism supports the preemption of an event
 4499 handler function by a higher priority event notification and supports the selection of the
 4500 highest priority pending event notification when multiple notifications (of different
 4501 priority) are pending simultaneously.

4502 The model here is based on hardware interrupts. Asynchronous event handling allows
 4503 the application to ensure that time-critical events are immediately processed when
 4504 delivered, without the indeterminism of being at a random location within a polling loop.

- 4505 Use of handler priority allows the specification of how handlers are interrupted by other
4506 higher priority handlers.
- 4507 — Differentiation between multiple occurrences of event notifications of the same type
- 4508 The events notification mechanism passes an application-defined value to the event
4509 handler function. This value can be used for a variety of purposes, such as enabling the
4510 application to identify which of several possible events of the same type (for example,
4511 timer expirations) has occurred.
- 4512 — Polled reception of asynchronous event notifications
- 4513 The events notification mechanism supports blocking and non-blocking polls for
4514 asynchronous event notification.
- 4515 The polled mode of operation is often preferred over the interrupt mode by those
4516 practitioners accustomed to this model. Providing support for this model facilitates the
4517 porting of applications based on this model to POSIX.1b conforming systems.
- 4518 — Deterministic response to asynchronous event notifications
- 4519 The events notification mechanism does not preclude implementations that provide
4520 deterministic event dispatch latency and minimizes the number of system calls needed to
4521 use the event facilities during realtime processing.
- 4522 • Rationale for Extension
- 4523 POSIX.1 signals have many of the characteristics necessary to support the asynchronous
4524 handling of event notifications, and the Realtime Signals Extension addresses the following
4525 deficiencies in the POSIX.1 signal mechanism:
- 4526 — Signals do not support reliable delivery of event notification. Subsequent occurrences of
4527 a pending signal are not guaranteed to be delivered.
- 4528 — Signals do not support prioritized delivery of event notifications. The order of signal
4529 delivery when multiple unblocked signals are pending is undefined.
- 4530 — Signals do not support the differentiation between multiple signals of the same type.

4531 *B.2.8.2 Asynchronous I/O*

4532 Many applications need to interact with the I/O subsystem in an asynchronous manner. The
4533 asynchronous I/O mechanism provides the ability to overlap application processing and I/O
4534 operations initiated by the application. The asynchronous I/O mechanism allows a single
4535 process to perform I/O simultaneously to a single file multiple times or to multiple files
4536 multiple times.

4537 **Overview**

4538 Asynchronous I/O operations proceed in logical parallel with the processing done by the
4539 application after the asynchronous I/O has been initiated. Other than this difference,
4540 asynchronous I/O behaves similarly to normal I/O using *read()*, *write()*, *lseek()*, and *fsync()*.
4541 The effect of issuing an asynchronous I/O request is as if a separate thread of execution were to
4542 perform atomically the implied *lseek()* operation, if any, and then the requested I/O operation
4543 (either *read()*, *write()*, or *fsync()*). There is no seek implied with a call to *aio_fsync()*. Concurrent
4544 asynchronous operations and synchronous operations applied to the same file update the file as
4545 if the I/O operations had proceeded serially.

4546 When asynchronous I/O completes, a signal can be delivered to the application to indicate the
4547 completion of the I/O. This signal can be used to indicate that buffers and control blocks used

4548 for asynchronous I/O can be reused. Signal delivery is not required for an asynchronous
 4549 operation and may be turned off on a per-operation basis by the application. Signals may also be
 4550 synchronously polled using *aio_suspend()*, *sigtimedwait()*, or *sigwaitinfo()*.

4551 Normal I/O has a return value and an error status associated with it. Asynchronous I/O returns
 4552 a value and an error status when the operation is first submitted, but that only relates to whether
 4553 the operation was successfully queued up for servicing. The I/O operation itself also has a
 4554 return status and an error value. To allow the application to retrieve the return status and the
 4555 error value, functions are provided that, given the address of an asynchronous I/O control
 4556 block, yield the return and error status associated with the operation. Until an asynchronous I/O
 4557 operation is done, its error status is [EINPROGRESS]. Thus, an application can poll for
 4558 completion of an asynchronous I/O operation by waiting for the error status to become equal to
 4559 a value other than [EINPROGRESS]. The return status of an asynchronous I/O operation is
 4560 undefined so long as the error status is equal to [EINPROGRESS].

4561 Storage for asynchronous operation return and error status may be limited. Submission of
 4562 asynchronous I/O operations may fail if this storage is exceeded. When an application retrieves
 4563 the return status of a given asynchronous operation, therefore, any system-maintained storage
 4564 used for this status and the error status may be reclaimed for use by other asynchronous
 4565 operations.

4566 Asynchronous I/O can be performed on file descriptors that have been enabled for POSIX.1b
 4567 synchronized I/O. In this case, the I/O operation still occurs asynchronously, as defined herein;
 4568 however, the asynchronous operation I/O in this case is not completed until the I/O has reached
 4569 either the state of synchronized I/O data integrity completion or synchronized I/O file integrity
 4570 completion, depending on the sort of synchronized I/O that is enabled on the file descriptor.

4571 Models

4572 Three models illustrate the use of asynchronous I/O: a journalization model, a data acquisition
 4573 model, and a model of the use of asynchronous I/O in supercomputing applications.

- 4574 • Journalization Model

4575 Many realtime applications perform low-priority journalizing functions. Journalizing
 4576 requires that logging records be queued for output without blocking the initiating process.

- 4577 • Data Acquisition Model

4578 A data acquisition process may also serve as a model. The process has two or more channels
 4579 delivering intermittent data that must be read within a certain time. The process issues one
 4580 asynchronous read on each channel. When one of the channels needs data collection, the
 4581 process reads the data and posts it through an asynchronous write to secondary memory for
 4582 future processing.

- 4583 • Supercomputing Model

4584 The supercomputing community has used asynchronous I/O much like that specified in
 4585 POSIX.1 for many years. This community requires the ability to perform multiple I/O
 4586 operations to multiple devices with a minimal number of entries to “the system”; each entry
 4587 to “the system” provokes a major delay in operations when compared to the normal progress
 4588 made by the application. This existing practice motivated the use of combined *lseek()* and
 4589 *read()* or *write()* calls, as well as the *lio_listio()* call. Another common practice is to disable
 4590 signal notification for I/O completion, and simply poll for I/O completion at some interval
 4591 by which the I/O should be completed. Likewise, interfaces like *aio_cancel()* have been in
 4592 successful commercial use for many years. Note also that an underlying implementation of
 4593 asynchronous I/O will require the ability, at least internally, to cancel outstanding

4594 asynchronous I/O, at least when the process exits. (Consider an asynchronous read from a
4595 terminal, when the process intends to exit immediately.)

4596 **Requirements**

4597 Asynchronous input and output for realtime implementations have these requirements:

- 4598 • The ability to queue multiple asynchronous read and write operations to a single open
4599 instance. Both sequential and random access should be supported.
- 4600 • The ability to queue asynchronous read and write operations to multiple open instances.
- 4601 • The ability to obtain completion status information by polling and/or asynchronous event
4602 notification.
- 4603 • Asynchronous event notification on asynchronous I/O completion is optional.
- 4604 • It has to be possible for the application to associate the event with the *aiocbp* for the operation
4605 that generated the event.
- 4606 • The ability to cancel queued requests.
- 4607 • The ability to wait upon asynchronous I/O completion in conjunction with other types of
4608 events.
- 4609 • The ability to accept an *aio_read()* and an *aio_cancel()* for a device that accepts a *read()*, and
4610 the ability to accept an *aio_write()* and an *aio_cancel()* for a device that accepts a *write()*. This
4611 does not imply that the operation is asynchronous.

4612 **Standardization Issues**

4613 The following issues are addressed by the standardization of asynchronous I/O:

4614 • Rationale for New Interface

4615 Non-blocking I/O does not satisfy the needs of either realtime or high-performance
4616 computing models; these models require that a process overlap program execution and I/O
4617 processing. Realtime applications will often make use of direct I/O to or from the address
4618 space of the process, or require synchronized (unbuffered) I/O; they also require the ability
4619 to overlap this I/O with other computation. In addition, asynchronous I/O allows an
4620 application to keep a device busy at all times, possibly achieving greater throughput.
4621 Supercomputing and database architectures will often have specialized hardware that can
4622 provide true asynchrony underlying the logical asynchrony provided by this interface. In
4623 addition, asynchronous I/O should be supported by all types of files and devices in the same
4624 manner.

4625 • Effect of Buffering

4626 If asynchronous I/O is performed on a file that is buffered prior to being actually written to
4627 the device, it is possible that asynchronous I/O will offer no performance advantage over
4628 normal I/O; the cycles *stolen* to perform the asynchronous I/O will be taken away from the
4629 running process and the I/O will occur at interrupt time. This potential lack of gain in
4630 performance in no way obviates the need for asynchronous I/O by realtime applications,
4631 which very often will use specialized hardware support, multiple processors, and/or
4632 unbuffered, synchronized I/O.

4633 **B.2.8.3 Memory Management**

4634 All memory management and shared memory definitions are located in the `<sys/mman.h>`
4635 header. This is for alignment with historical practice.

4636 **Memory Locking Functions**

4637 This portion of the rationale presents models, requirements, and standardization issues relevant
4638 to process memory locking.

4639 • Models

4640 Realtime systems that conform to IEEE Std 1003.1-2001 are expected (and desired) to be
4641 supported on systems with demand-paged virtual memory management, non-paged
4642 swapping memory management, and physical memory systems with no memory
4643 management hardware. The general case, however, is the demand-paged, virtual memory
4644 system with each POSIX process running in a virtual address space. Note that this includes
4645 architectures where each process resides in its own virtual address space and architectures
4646 where the address space of each process is only a portion of a larger global virtual address
4647 space.

4648 The concept of memory locking is introduced to eliminate the indeterminacy introduced by
4649 paging and swapping, and to support an upper bound on the time required to access the
4650 memory mapped into the address space of a process. Ideally, this upper bound will be the
4651 same as the time required for the processor to access “main memory”, including any address
4652 translation and cache miss overheads. But some implementations—primarily on
4653 mainframes—will not actually force locked pages to be loaded and held resident in main
4654 memory. Rather, they will handle locked pages so that accesses to these pages will meet the
4655 performance metrics for locked process memory in the implementation. Also, although it is
4656 not, for example, the intention that this interface, as specified, be used to lock process
4657 memory into “cache”, it is conceivable that an implementation could support a large static
4658 RAM memory and define this as “main memory” and use a large[r] dynamic RAM as
4659 “backing store”. These interfaces could then be interpreted as supporting the locking of
4660 process memory into the static RAM. Support for multiple levels of backing store would
4661 require extensions to these interfaces.

4662 Implementations may also use memory locking to guarantee a fixed translation between
4663 virtual and physical addresses where such is beneficial to improving determinacy for
4664 direct-to/from-process input/output. IEEE Std 1003.1-2001 does not guarantee to the
4665 application that the virtual-to-physical address translations, if such exist, are fixed, because
4666 such behavior would not be implementable on all architectures on which implementations of
4667 IEEE Std 1003.1-2001 are expected. But IEEE Std 1003.1-2001 does mandate that an
4668 implementation define, for the benefit of potential users, whether or not locking guarantees
4669 fixed translations.

4670 Memory locking is defined with respect to the address space of a process. Only the pages
4671 mapped into the address space of a process may be locked by the process, and when the
4672 pages are no longer mapped into the address space—for whatever reason—the locks
4673 established with respect to that address space are removed. Shared memory areas warrant
4674 special mention, as they may be mapped into more than one address space or mapped more
4675 than once into the address space of a process; locks may be established on pages within these
4676 areas with respect to several of these mappings. In such a case, the lock state of the
4677 underlying physical pages is the logical OR of the lock state with respect to each of the
4678 mappings. Only when all such locks have been removed are the shared pages considered
4679 unlocked.

4680 In recognition of the page granularity of Memory Management Units (MMU), and in order to
4681 support locking of ranges of address space, memory locking is defined in terms of “page”
4682 granularity. That is, for the interfaces that support an address and size specification for the
4683 region to be locked, the address must be on a page boundary, and all pages mapped by the
4684 specified range are locked, if valid. This means that the length is implicitly rounded up to a
4685 multiple of the page size. The page size is implementation-defined and is available to
4686 applications as a compile-time symbolic constant or at runtime via *sysconf()*.

4687 A “real memory” POSIX.1b implementation that has no MMU could elect not to support
4688 these interfaces, returning [ENOSYS]. But an application could easily interpret this as
4689 meaning that the implementation would unconditionally page or swap the application when
4690 such is not the case. It is the intention of IEEE Std 1003.1-2001 that such a system could define
4691 these interfaces as “NO-OPs”, returning success without actually performing any function
4692 except for mandated argument checking.

4693 • Requirements

4694 For realtime applications, memory locking is generally considered to be required as part of
4695 application initialization. This locking is performed after an application has been loaded (that
4696 is, *exec'd*) and the program remains locked for its entire lifetime. But to support applications
4697 that undergo major mode changes where, in one mode, locking is required, but in another it
4698 is not, the specified interfaces allow repeated locking and unlocking of memory within the
4699 lifetime of a process.

4700 When a realtime application locks its address space, it should not be necessary for the
4701 application to then “touch” all of the pages in the address space to guarantee that they are
4702 resident or else suffer potential paging delays the first time the page is referenced. Thus,
4703 IEEE Std 1003.1-2001 requires that the pages locked by the specified interfaces be resident
4704 when the locking functions return successfully.

4705 Many architectures support system-managed stacks that grow automatically when the
4706 current extent of the stack is exceeded. A realtime application has a requirement to be able to
4707 “preallocate” sufficient stack space and lock it down so that it will not suffer page faults to
4708 grow the stack during critical realtime operation. There was no consensus on a portable way
4709 to specify how much stack space is needed, so IEEE Std 1003.1-2001 supports no specific
4710 interface for preallocating stack space. But an application can portably lock down a specific
4711 amount of stack space by specifying *MCL_FUTURE* in a call to *mlockall()* and then calling a
4712 dummy function that declares an automatic array of the desired size.

4713 Memory locking for realtime applications is also generally considered to be an “all or
4714 nothing” proposition. That is, the entire process, or none, is locked down. But, for
4715 applications that have well-defined sections that need to be locked and others that do not,
4716 IEEE Std 1003.1-2001 supports an optional set of interfaces to lock or unlock a range of
4717 process addresses. Reasons for locking down a specific range include:

4718 — An asynchronous event handler function that must respond to external events in a
4719 deterministic manner such that page faults cannot be tolerated

4720 — An input/output “buffer” area that is the target for direct-to-process I/O, and the
4721 overhead of implicit locking and unlocking for each I/O call cannot be tolerated

4722 Finally, locking is generally viewed as an “application-wide” function. That is, the
4723 application is globally aware of which regions are locked and which are not over time. This is
4724 in contrast to a function that is used temporarily within a “third party” library routine whose
4725 function is unknown to the application, and therefore must have no “side effects”. The
4726 specified interfaces, therefore, do not support “lock stacking” or “lock nesting” within a
4727 process. But, for pages that are shared between processes or mapped more than once into a

4728 process address space, “lock stacking” is essentially mandated by the requirement that
4729 unlocking of pages that are mapped by more than one process or more than once by the same
4730 process does not affect locks established on the other mappings.

4731 There was some support for “lock stacking” so that locking could be transparently used in
4732 functions or opaque modules. But the consensus was not to burden all implementations with
4733 lock stacking (and reference counting), and an implementation option was proposed. There
4734 were strong objections to the option because applications would have to support both
4735 options in order to remain portable. The consensus was to eliminate lock stacking altogether,
4736 primarily through overwhelming support for the System V “m[un]lock[all]” interface on
4737 which IEEE Std 1003.1-2001 is now based.

4738 Locks are not inherited across *fork()*s because some implementations implement *fork()* by
4739 creating new address spaces for the child. In such an implementation, requiring locks to be
4740 inherited would lead to new situations in which a fork would fail due to the inability of the
4741 system to lock sufficient memory to lock both the parent and the child. The consensus was
4742 that there was no benefit to such inheritance. Note that this does not mean that locks are
4743 removed when, for instance, a thread is created in the same address space.

4744 Similarly, locks are not inherited across *exec* because some implementations implement *exec*
4745 by unmapping all of the pages in the address space (which, by definition, removes the locks
4746 on these pages), and maps in pages of the *exec*'d image. In such an implementation, requiring
4747 locks to be inherited would lead to new situations in which *exec* would fail. Reporting this
4748 failure would be very cumbersome to detect in time to report to the calling process, and no
4749 appropriate mechanism exists for informing the *exec*'d process of its status.

4750 It was determined that, if the newly loaded application required locking, it was the
4751 responsibility of that application to establish the locks. This is also in keeping with the
4752 general view that it is the responsibility of the application to be aware of all locks that are
4753 established.

4754 There was one request to allow (not mandate) locks to be inherited across *fork()*, and a
4755 request for a flag, MCL_INHERIT, that would specify inheritance of memory locks across
4756 *exec*s. Given the difficulties raised by this and the general lack of support for the feature in
4757 IEEE Std 1003.1-2001, it was not added. IEEE Std 1003.1-2001 does not preclude an
4758 implementation from providing this feature for administrative purposes, such as a “run”
4759 command that will lock down and execute a specified application. Additionally, the rationale
4760 for the objection equated *fork()* with creating a thread in the address space.
4761 IEEE Std 1003.1-2001 does not mandate releasing locks when creating additional threads in
4762 an existing process.

4763 • Standardization Issues

4764 One goal of IEEE Std 1003.1-2001 is to define a set of primitives that provide the necessary
4765 functionality for realtime applications, with consideration for the needs of other application
4766 domains where such were identified, which is based to the extent possible on existing
4767 industry practice.

4768 The Memory Locking option is required by many realtime applications to tune performance.
4769 Such a facility is accomplished by placing constraints on the virtual memory system to limit
4770 paging of time of the process or of critical sections of the process. This facility should not be
4771 used by most non-realtime applications.

4772 Optional features provided in IEEE Std 1003.1-2001 allow applications to lock selected
4773 address ranges with the caveat that the process is responsible for being aware of the page
4774 granularity of locking and the unnested nature of the locks.

4775 **Mapped Files Functions**

4776 The Memory Mapped Files option provides a mechanism that allows a process to access files by
4777 directly incorporating file data into its address space. Once a file is “mapped” into a process
4778 address space, the data can be manipulated by instructions as memory. The use of mapped files
4779 can significantly reduce I/O data movement since file data does not have to be copied into
4780 process data buffers as in *read()* and *write()*. If more than one process maps a file, its contents
4781 are shared among them. This provides a low overhead mechanism by which processes can
4782 synchronize and communicate.

4783 • Historical Perspective

4784 Realtime applications have historically been implemented using a collection of cooperating
4785 processes or tasks. In early systems, these processes ran on bare hardware (that is, without an
4786 operating system) with no memory relocation or protection. The application paradigms that
4787 arose from this environment involve the sharing of data between the processes.

4788 When realtime systems were implemented on top of vendor-supplied operating systems, the
4789 paradigm or performance benefits of direct access to data by multiple processes was still
4790 deemed necessary. As a result, operating systems that claim to support realtime applications
4791 must support the shared memory paradigm.

4792 Additionally, a number of realtime systems provide the ability to map specific sections of the
4793 physical address space into the address space of a process. This ability is required if an
4794 application is to obtain direct access to memory locations that have specific properties (for
4795 example, refresh buffers or display devices, dual ported memory locations, DMA target
4796 locations). The use of this ability is common enough to warrant some degree of
4797 standardization of its interface. This ability overlaps the general paradigm of shared
4798 memory in that, in both instances, common global objects are made addressable by
4799 individual processes or tasks.

4800 Finally, a number of systems also provide the ability to map process addresses to files. This
4801 provides both a general means of sharing persistent objects, and using files in a manner that
4802 optimizes memory and swapping space usage.

4803 Simple shared memory is clearly a special case of the more general file mapping capability.
4804 In addition, there is relatively widespread agreement and implementation of the file
4805 mapping interface. In these systems, many different types of objects can be mapped (for
4806 example, files, memory, devices, and so on) using the same mapping interfaces. This
4807 approach both minimizes interface proliferation and maximizes the generality of programs
4808 using the mapping interfaces.

4809 • Memory Mapped Files Usage

4810 A memory object can be concurrently mapped into the address space of one or more
4811 processes. The *mmap()* and *munmap()* functions allow a process to manipulate their address
4812 space by mapping portions of memory objects into it and removing them from it. When
4813 multiple processes map the same memory object, they can share access to the underlying
4814 data. Implementations may restrict the size and alignment of mappings to be on *page-size*
4815 boundaries. The page size, in bytes, is the value of the system-configurable variable
4816 {PAGESIZE}, typically accessed by calling *sysconf()* with a *name* argument of
4817 *_SC_PAGESIZE*. If an implementation has no restrictions on size or alignment, it may
4818 specify a 1-byte page size.

4819 To map memory, a process first opens a memory object. The *ftruncate()* function can be used
4820 to contract or extend the size of the memory object even when the object is currently
4821 mapped. If the memory object is extended, the contents of the extended areas are zeros.

4822 After opening a memory object, the application maps the object into its address space using
4823 the *mmap()* function call. Once a mapping has been established, it remains mapped until
4824 unmapped with *munmap()*, even if the memory object is closed. The *mprotect()* function can
4825 be used to change the memory protections initially established by *mmap()*.

4826 A *close()* of the file descriptor, while invalidating the file descriptor itself, does not unmap
4827 any mappings established for the memory object. The address space, including all mapped
4828 regions, is inherited on *fork()*. The entire address space is unmapped on process termination
4829 or by successful calls to any of the *exec* family of functions.

4830 The *msync()* function is used to force mapped file data to permanent storage.

4831 • Effects on Other Functions

4832 When the Memory Mapped Files option is supported, the operation of the *open()*, *creat()*, and
4833 *unlink()* functions are a natural result of using the file system name space to map the global
4834 names for memory objects.

4835 The *ftruncate()* function can be used to set the length of a sharable memory object.

4836 The meaning of *stat()* fields other than the size and protection information is undefined on
4837 implementations where memory objects are not implemented using regular files. When
4838 regular files are used, the times reflect when the implementation updated the file image of
4839 the data, not when a process updated the data in memory.

4840 The operations of *fdopen()*, *write()*, *read()*, and *lseek()* were made unspecified for objects
4841 opened with *shm_open()*, so that implementations that did not implement memory objects as
4842 regular files would not have to support the operation of these functions on shared memory
4843 objects.

4844 The behavior of memory objects with respect to *close()*, *dup()*, *dup2()*, *open()*, *close()*, *fork()*,
4845 *_exit()*, and the *exec* family of functions is the same as the behavior of the existing practice of
4846 the *mmap()* function.

4847 A memory object can still be referenced after a close. That is, any mappings made to the file
4848 are still in effect, and reads and writes that are made to those mappings are still valid and are
4849 shared with other processes that have the same mapping. Likewise, the memory object can
4850 still be used if any references remain after its name(s) have been deleted. Any references that
4851 remain after a close must not appear to the application as file descriptors.

4852 This is existing practice for *mmap()* and *close()*. In addition, there are already mappings
4853 present (text, data, stack) that do not have open file descriptors. The text mapping in
4854 particular is considered a reference to the file containing the text. The desire was to treat all
4855 mappings by the process uniformly. Also, many modern implementations use *mmap()* to
4856 implement shared libraries, and it would not be desirable to keep file descriptors for each of
4857 the many libraries an application can use. It was felt there were many other existing
4858 programs that used this behavior to free a file descriptor, and thus IEEE Std 1003.1-2001
4859 could not forbid it and still claim to be using existing practice.

4860 For implementations that implement memory objects using memory only, memory objects
4861 will retain the memory allocated to the file after the last close and will use that same memory
4862 on the next open. Note that closing the memory object is not the same as deleting the name,
4863 since the memory object is still defined in the memory object name space.

4864 The locks of *fcntl()* do not block any read or write operation, including read or write access to
4865 shared memory or mapped files. In addition, implementations that only support shared
4866 memory objects should not be required to implement record locks. The reference to *fcntl()* is
4867 added to make this point explicitly. The other *fcntl()* commands are useful with shared

4868 memory objects.

4869 The size of pages that mapping hardware may be able to support may be a configurable
4870 value, or it may change based on hardware implementations. The addition of the
4871 `_SC_PAGESIZE` parameter to the `sysconf()` function is provided for determining the mapping
4872 page size at runtime.

4873 Shared Memory Functions

4874 Implementations may support the Shared Memory Objects option without supporting a general
4875 Memory Mapped Files option. Shared memory objects are named regions of storage that may be
4876 independent of the file system and can be mapped into the address space of one or more
4877 processes to allow them to share the associated memory.

4878 • Requirements

4879 Shared memory is used to share data among several processes, each potentially running at
4880 different priority levels, responding to different inputs, or performing separate tasks. Shared
4881 memory is not just simply providing common access to data, it is providing the fastest
4882 possible communication between the processes. With one memory write operation, a process
4883 can pass information to as many processes as have the memory region mapped.

4884 As a result, shared memory provides a mechanism that can be used for all other interprocess
4885 communication facilities. It may also be used by an application for implementing more
4886 sophisticated mechanisms than semaphores and message queues.

4887 The need for a shared memory interface is obvious for virtual memory systems, where the
4888 operating system is directly preventing processes from accessing each other's data. However,
4889 in unprotected systems, such as those found in some embedded controllers, a shared
4890 memory interface is needed to provide a portable mechanism to allocate a region of memory
4891 to be shared and then to communicate the address of that region to other processes.

4892 This, then, provides the minimum functionality that a shared memory interface must have in
4893 order to support realtime applications: to allocate and name an object to be mapped into
4894 memory for potential sharing (`open()` or `shm_open()`), and to make the memory object
4895 available within the address space of a process (`mmap()`). To complete the interface, a
4896 mechanism to release the claim of a process on a shared memory object (`munmap()`) is also
4897 needed, as well as a mechanism for deleting the name of a sharable object that was
4898 previously created (`unlink()` or `shm_unlink()`).

4899 After a mapping has been established, an implementation should not have to provide
4900 services to maintain that mapping. All memory writes into that area will appear immediately
4901 in the memory mapping of that region by any other processes.

4902 Thus, requirements include:

- 4903 — Support creation of sharable memory objects and the mapping of these objects into the
4904 address space of a process.
- 4905 — Sharable memory objects should be accessed by global names accessible from all
4906 processes.
- 4907 — Support the mapping of specific sections of physical address space (such as a memory
4908 mapped device) into the address space of a process. This should not be done by the
4909 process specifying the actual address, but again by an implementation-defined global
4910 name (such as a special device name) dedicated to this purpose.
- 4911 — Support the mapping of discrete portions of these memory objects.

- 4912 — Support for minimum hardware configurations that contain no physical media on which
4913 to store shared memory contents permanently.
- 4914 — The ability to preallocate the entire shared memory region so that minimum hardware
4915 configurations without virtual memory support can guarantee contiguous space.
- 4916 — The maximizing of performance by not requiring functionality that would require
4917 implementation interaction above creating the shared memory area and returning the
4918 mapping.
- 4919 Note that the above requirements do not preclude:
- 4920 — The sharable memory object from being implemented using actual files on an actual file
4921 system.
- 4922 — The global name that is accessible from all processes being restricted to a file system area
4923 that is dedicated to handling shared memory.
- 4924 — An implementation not providing implementation-defined global names for the purpose
4925 of physical address mapping.
- 4926 • Shared Memory Objects Usage
- 4927 If the Shared Memory Objects option is supported, a shared memory object may be created,
4928 or opened if it already exists, with the *shm_open()* function. If the shared memory object is
4929 created, it has a length of zero. The *truncate()* function can be used to set the size of the
4930 shared memory object after creation. The *shm_unlink()* function removes the name for a
4931 shared memory object created by *shm_open()*.
- 4932 • Shared Memory Overview
- 4933 The shared memory facility defined by IEEE Std 1003.1-2001 usually results in memory
4934 locations being added to the address space of the process. The implementation returns the
4935 address of the new space to the application by means of a pointer. This works well in
4936 languages like C. However, in languages without pointer types it will not work. In the
4937 bindings for such a language, either a special COMMON section will need to be defined
4938 (which is unlikely), or the binding will have to allow existing structures to be mapped. The
4939 implementation will likely have to place restrictions on the size and alignment of such
4940 structures or will have to map a suitable region of the address space of the process into the
4941 memory object, and thus into other processes. These are issues for that particular language
4942 binding. For IEEE Std 1003.1-2001, however, the practice will not be forbidden, merely
4943 undefined.
- 4944 Two potentially different name spaces are used for naming objects that may be mapped into
4945 process address spaces. When the Memory Mapped Files option is supported, files may be
4946 accessed via *open()*. When the Shared Memory Objects option is supported, sharable
4947 memory objects that might not be files may be accessed via the *shm_open()* function. These
4948 options are not mutually-exclusive.
- 4949 Some implementations supporting the Shared Memory Objects option may choose to
4950 implement the shared memory object name space as part of the file system name space.
4951 There are several reasons for this:
- 4952 — It allows applications to prevent name conflicts by use of the directory structure.
- 4953 — It uses an existing mechanism for accessing global objects and prevents the creation of a
4954 new mechanism for naming global objects.
- 4955 In such implementations, memory objects can be implemented using regular files, if that is
4956 what the implementation chooses. The *shm_open()* function can be implemented as an *open()*

4957 call in a fixed directory followed by a call to *fcntl()* to set *FD_CLOEXEC*. The *shm_unlink()*
4958 function can be implemented as an *unlink()* call.

4959 On the other hand, it is also expected that small embedded systems that support the Shared
4960 Memory Objects option may wish to implement shared memory without having any file
4961 systems present. In this case, the implementations may choose to use a simple string valued
4962 name space for shared memory regions. The *shm_open()* function permits either type of
4963 implementation.

4964 Some implementations have hardware that supports protection of mapped data from certain
4965 classes of access and some do not. Systems that supply this functionality can support the
4966 Memory Protection option.

4967 Some implementations restrict size, alignment, and protections to be on *page*-size
4968 boundaries. If an implementation has no restrictions on size or alignment, it may specify a 1-
4969 byte page size. Applications on implementations that do support larger pages must be
4970 cognizant of the page size since this is the alignment and protection boundary.

4971 Simple embedded implementations may have a 1-byte page size and only support the Shared
4972 Memory Objects option. This provides simple shared memory between processes without
4973 requiring mapping hardware.

4974 IEEE Std 1003.1-2001 specifically allows a memory object to remain referenced after a close
4975 because that is existing practice for the *mmap()* function.

4976 **Typed Memory Functions**

4977 Implementations may support the Typed Memory Objects option without supporting either the
4978 Shared Memory option or the Memory Mapped Files option. Typed memory objects are pools of
4979 specialized storage, different from the main memory resource normally used by a processor to
4980 hold code and data, that can be mapped into the address space of one or more processes.

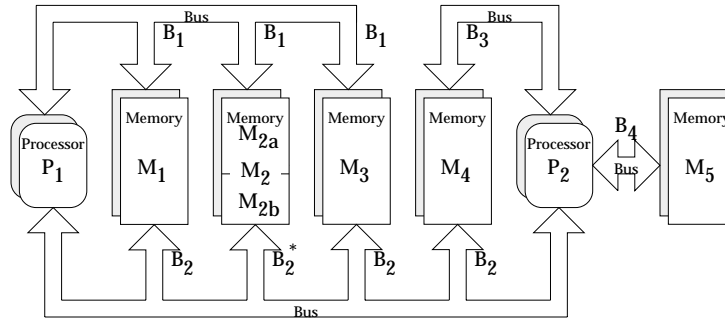
4981 • Model

4982 Realtime systems conforming to one of the POSIX.13 realtime profiles are expected (and
4983 desired) to be supported on systems with more than one type or pool of memory (for
4984 example, SRAM, DRAM, ROM, EPROM, EEPROM), where each type or pool of memory may
4985 be accessible by one or more processors via one or more busses (ports). Memory mapped
4986 files, shared memory objects, and the language-specific storage allocation operators (*malloc()*
4987 for the ISO C standard, *new* for ISO Ada) fail to provide application program interfaces
4988 versatile enough to allow applications to control their utilization of such diverse memory
4989 resources. The typed memory interfaces *posix_typed_mem_open()*, *posix_mem_offset()*,
4990 *posix_typed_mem_get_info()*, *mmap()*, and *munmap()* defined herein support the model of
4991 typed memory described below.

4992 For purposes of this model, a system comprises several processors (for example, P_1 and P_2),
4993 several physical memory pools (for example, M_1 , M_2 , M_{2a} , M_{2b} , M_3 , M_4 , and M_5), and several
4994 busses or “ports” (for example, B_1 , B_2 , B_3 , and B_4) interconnecting the various processors and
4995 memory pools in some system-specific way. Notice that some memory pools may be
4996 contained in others (for example, M_{2a} and M_{2b} are contained in M_2).

4997 Figure B-1 (on page 124) shows an example of such a model. In a system like this, an
4998 application should be able to perform the following operations:

4999



* All addresses in pool M_2 (comprising pools M_{2a} and M_{2b}) accessible via port B_1 .
 Addresses in pool M_{2b} are also accessible via port B_2 .
 Addresses in pool M_{2a} are *not* accessible via port B_2 .

5000

Figure B-1 Example of a System with Typed Memory

5001

— Typed Memory Allocation

5002

5003

5004

5005

5006

5007

An application should be able to allocate memory dynamically from the desired pool using the desired bus, and map it into a process' address space. For example, processor P_1 can allocate some portion of memory pool M_1 through port B_1 , treating all unmapped subareas of M_1 as a heap-storage resource from which memory may be allocated. This portion of memory is mapped into the process' address space, and subsequently deallocated when unmapped from all processes.

5008

— Using the Same Storage Region from Different Busses

5009

5010

5011

5012

5013

An application process with a mapped region of storage that is accessed from one bus should be able to map that same storage area at another address (subject to page size restrictions detailed in *mmmap()*), to allow it to be accessed from another bus. For example, processor P_1 may wish to access the same region of memory pool M_{2b} both through ports B_1 and B_2 .

5014

— Sharing Typed Memory Regions

5015

5016

5017

5018

5019

5020

5021

5022

5023

5024

5025

5026

Several application processes running on the same or different processors may wish to share a particular region of a typed memory pool. Each process or processor may wish to access this region through different busses. For example, processor P_1 may want to share a region of memory pool M_4 with processor P_2 , and they may be required to use busses B_2 and B_3 , respectively, to minimize bus contention. A problem arises here when a process allocates and maps a portion of fragmented memory and then wants to share this region of memory with another process, either in the same processor or different processors. The solution adopted is to allow the first process to find out the memory map (offsets and lengths) of all the different fragments of memory that were mapped into its address space, by repeatedly calling *posix_mem_offset()*. Then, this process can pass the offsets and lengths obtained to the second process, which can then map the same memory fragments into its address space.

5027

— Contiguous Allocation

5028

5029

5030

5031

The problem of finding the memory map of the different fragments of the memory pool that were mapped into logically contiguous addresses of a given process can be solved by requesting contiguous allocation. For example, a process in P_1 can allocate 10 Kbytes of physically contiguous memory from M_3 - B_1 , and obtain the offset (within pool M_3) of this

5032 block of memory. Then, it can pass this offset (and the length) to a process in P_2 using
 5033 some interprocess communication mechanism. The second process can map the same
 5034 block of memory by using the offset transferred and specifying M_3-B_2 .

5035 — Unallocated Mapping

5036 Any subarea of a memory pool that is mapped to a process, either as the result of an
 5037 allocation request or an explicit mapping, is normally unavailable for allocation. Special
 5038 processes such as debuggers, however, may need to map large areas of a typed memory
 5039 pool, yet leave those areas available for allocation.

5040 Typed memory allocation and mapping has to coexist with storage allocation operators like
 5041 *malloc()*, but systems are free to choose how to implement this coexistence. For example, it
 5042 may be system configuration-dependent if all available system memory is made part of one
 5043 of the typed memory pools or if some part will be restricted to conventional allocation
 5044 operators. Equally system configuration-dependent may be the availability of operators like
 5045 *malloc()* to allocate storage from certain typed memory pools. It is not excluded to configure
 5046 a system such that a given named pool, P_1 , is in turn split into non-overlapping named
 5047 subpools. For example, M_1-B_1 , M_2-B_1 , and M_3-B_1 could also be accessed as one common pool
 5048 $M_{123}-B_1$. A call to *malloc()* on P_1 could work on such a larger pool while full optimization of
 5049 memory usage by P_1 would require typed memory allocation at the subpool level.

5050 • Existing Practice

5051 OS-9 provides for the naming (numbering) and prioritization of memory types by a system
 5052 administrator. It then provides APIs to request memory allocation of typed (colored)
 5053 memory by number, and to generate a bus address from a mapped memory address
 5054 (translate). When requesting colored memory, the user can specify type 0 to signify allocation
 5055 from the first available type in priority order.

5056 HP-RT presents interfaces to map different kinds of storage regions that are visible through a
 5057 VME bus, although it does not provide allocation operations. It also provides functions to
 5058 perform address translation between VME addresses and virtual addresses. It represents a
 5059 VME-bus unique solution to the general problem.

5060 The PSOS approach is similar (that is, based on a pre-established mapping of bus address
 5061 ranges to specific memories) with a concept of segments and regions (regions dynamically
 5062 allocated from a heap which is a special segment). Therefore, PSOS does not fully address the
 5063 general allocation problem either. PSOS does not have a “process”-based model, but more of
 5064 a “thread”-only-based model of multi-tasking. So mapping to a process address space is not
 5065 an issue.

5066 QNX uses the System V approach of opening specially named devices (shared memory
 5067 segments) and using *mmap()* to then gain access from the process. They do not address
 5068 allocation directly, but once typed shared memory can be mapped, an “allocation manager”
 5069 process could be written to handle requests for allocation.

5070 The System V approach also included allocation, implemented by opening yet other special
 5071 “devices” which allocate, rather than appearing as a whole memory object.

5072 The Orkid realtime kernel interface definition has operations to manage memory “regions”
 5073 and “pools”, which are areas of memory that may reflect the differing physical nature of the
 5074 memory. Operations to allocate memory from these regions and pools are also provided.

- 5075 • Requirements
- 5076 Existing practice in SVID-derived UNIX systems relies on functionality similar to *mmap()*
- 5077 and its related interfaces to achieve mapping and allocation of typed memory. However, the
- 5078 issue of sharing typed memory (allocated or mapped) and the complication of multiple ports
- 5079 are not addressed in any consistent way by existing UNIX system practice. Part of this
- 5080 functionality is existing practice in specialized realtime operating systems. In order to
- 5081 solidify the capabilities implied by the model above, the following requirements are imposed
- 5082 on the interface:
- 5083 — Identification of Typed Memory Pools and Ports
- 5084 All processes (running in all processors) in the system are able to identify a particular
- 5085 (system configured) typed memory pool accessed through a particular (system
- 5086 configured) port by a name. That name is a member of a name space common to all these
- 5087 processes, but need not be the same name space as that containing ordinary filenames.
- 5088 The association between memory pools/ports and corresponding names is typically
- 5089 established when the system is configured. The “open” operation for typed memory
- 5090 objects should be distinct from the *open()* function, for consistency with other similar
- 5091 services, but implementable on top of *open()*. This implies that the handle for a typed
- 5092 memory object will be a file descriptor.
- 5093 — Allocation and Mapping of Typed Memory
- 5094 Once a typed memory object has been identified by a process, it is possible to both map
- 5095 user-selected subareas of that object into process address space and to map system-
- 5096 selected (that is, dynamically allocated) subareas of that object, with user-specified
- 5097 length, into process address space. It is also possible to determine the maximum length of
- 5098 memory allocation that may be requested from a given typed memory object.
- 5099 — Sharing Typed Memory
- 5100 Two or more processes are able to share portions of typed memory, either user-selected or
- 5101 dynamically allocated. This requirement applies also to dynamically allocated regions of
- 5102 memory that are composed of several non-contiguous pieces.
- 5103 — Contiguous Allocation
- 5104 For dynamic allocation, it is the user’s option whether the system is required to allocate a
- 5105 contiguous subarea within the typed memory object, or whether it is permitted to allocate
- 5106 discontiguous fragments which appear contiguous in the process mapping. Contiguous
- 5107 allocation simplifies the process of sharing allocated typed memory, while discontiguous
- 5108 allocation allows for potentially better recovery of deallocated typed memory.
- 5109 — Accessing Typed Memory Through Different Ports
- 5110 Once a subarea of a typed memory object has been mapped, it is possible to determine the
- 5111 location and length corresponding to a user-selected portion of that object within the
- 5112 memory pool. This location and length can then be used to remap that portion of memory
- 5113 for access from another port. If the referenced portion of typed memory was allocated
- 5114 discontiguously, the length thus determined may be shorter than anticipated, and the
- 5115 user code must adapt to the value returned.
- 5116 — Deallocation
- 5117 When a previously mapped subarea of typed memory is no longer mapped by any
- 5118 process in the system—as a result of a call or calls to *munmap()*—that subarea becomes
- 5119 potentially reusable for dynamic allocation; actual reuse of the subarea is a function of the
- 5120 dynamic typed memory allocation policy.

5121 — Unallocated Mapping

5122 It must be possible to map user-selected subareas of a typed memory object without
 5123 marking that subarea as unavailable for allocation. This option is not the default behavior,
 5124 and requires appropriate privilege.

5125 • Scenario

5126 The following scenario will serve to clarify the use of the typed memory interfaces.

5127 Process A running on P_1 (see Figure B-1 (on page 124)) wants to allocate some memory from
 5128 memory pool M_2 , and it wants to share this portion of memory with process B running on P_2 .
 5129 Since P_2 only has access to the lower part of M_2 , both processes will use the memory pool
 5130 named M_{2b} which is the part of M_2 that is accessible both from P_1 and P_2 . The operations that
 5131 both processes need to perform are shown below:

5132 — Allocating Typed Memory

5133 Process A calls *posix_typed_mem_open()* with the name **/typed.m2b-b1** and a *tflag* of
 5134 POSIX_TYPED_MEM_ALLOCATE to get a file descriptor usable for allocating from pool
 5135 M_{2b} accessed through port B_1 . It then calls *mmap()* with this file descriptor requesting a
 5136 length of 4 096 bytes. The system allocates two discontinuous blocks of sizes 1 024 and
 5137 3 072 bytes within M_{2b} . The *mmap()* function returns a pointer to a 4 096-byte array in
 5138 process A's logical address space, mapping the allocated blocks contiguously. Process A
 5139 can then utilize the array, and store data in it.

5140 — Determining the Location of the Allocated Blocks

5141 Process A can determine the lengths and offsets (relative to M_{2b}) of the two blocks
 5142 allocated, by using the following procedure: First, process A calls *posix_mem_offset()* with
 5143 the address of the first element of the array and length 4 096. Upon return, the offset and
 5144 length (1 024 bytes) of the first block are returned. A second call to *posix_mem_offset()* is
 5145 then made using the address of the first element of the array plus 1 024 (the length of the
 5146 first block), and a new length of 4 096–1 024. If there were more fragments allocated, this
 5147 procedure could have been continued within a loop until the offsets and lengths of all the
 5148 blocks were obtained. Notice that this relatively complex procedure can be avoided if
 5149 contiguous allocation is requested (by opening the typed memory object with the *tflag*
 5150 POSIX_TYPED_MEM_ALLOCATE_CONTIG).

5151 — Sharing Data Across Processes

5152 Process A passes the two offset values and lengths obtained from the *posix_mem_offset()*
 5153 calls to process B running on P_2 , via some form of interprocess communication. Process B
 5154 can gain access to process A's data by calling *posix_typed_mem_open()* with the name
 5155 **/typed.m2b-b2** and a *tflag* of zero, then using two *mmap()* calls on the resulting file
 5156 descriptor to map the two subareas of that typed memory object to its own address space.

5157 • Rationale for no *mem_alloc()* and *mem_free()*

5158 The standard developers had originally proposed a pair of new flags to *mmap()* which, when
 5159 applied to a typed memory object descriptor, would cause *mmap()* to allocate dynamically
 5160 from an unallocated and unmapped area of the typed memory object. Deallocation was
 5161 similarly accomplished through the use of *munmap()*. This was rejected by the ballot group
 5162 because it excessively complicated the (already rather complex) *mmap()* interface and
 5163 introduced semantics useful only for typed memory, to a function which must also map
 5164 shared memory and files. They felt that a memory allocator should be built on top of *mmap()*
 5165 instead of being incorporated within the same interface, much as the ISO C standard libraries
 5166 build *malloc()* on top of the virtual memory mapping functions *brk()* and *sbrk()*. This would

5167 eliminate the complicated semantics involved with unmapping only part of an allocated
5168 block of typed memory.

5169 To attempt to achieve ballot group consensus, typed memory allocation and deallocation was
5170 first migrated from *mmap()* and *munmap()* to a pair of complementary functions modeled on
5171 the ISO C standard *malloc()* and *free()*. The *mem_alloc()* function specified explicitly the
5172 typed memory object (typed memory pool/access port) from which allocation takes place,
5173 unlike *malloc()* where the memory pool and port are unspecified. The *mem_free()* function
5174 handled deallocation. These new semantics still met all of the requirements detailed above
5175 without modifying the behavior of *mmap()* except to allow it to map specified areas of typed
5176 memory objects. An implementation would have been free to implement *mem_alloc()* and
5177 *mem_free()* over *mmap()*, through *mmap()*, or independently but cooperating with *mmap()*.

5178 The ballot group was queried to see if this was an acceptable alternative, and while there was
5179 some agreement that it achieved the goal of removing the complicated semantics of
5180 allocation from the *mmap()* interface, several balloters realized that it just created two
5181 additional functions that behaved, in great part, like *mmap()*. These balloters proposed an
5182 alternative which has been implemented here in place of a separate *mem_alloc()* and
5183 *mem_free()*. This alternative is based on four specific suggestions:

- 5184 1. The *posix_typed_mem_open()* function should provide a flag which specifies “allocate
5185 on *mmap()*” (otherwise, *mmap()* just maps the underlying object). This allows things
5186 roughly similar to */dev/zero* versus */dev/swap*. Two such flags have been implemented,
5187 one of which forces contiguous allocation.
- 5188 2. The *posix_mem_offset()* function is acceptable because it can be applied usefully to
5189 mapped objects in general. It should return the file descriptor of the underlying object.
- 5190 3. The *mem_get_info()* function in an earlier draft should be renamed
5191 *posix_typed_mem_get_info()* because it is not generally applicable to memory objects. It
5192 should probably return the file descriptor’s allocation attribute. The renaming of the
5193 function has been implemented, but having it return a piece of information which is
5194 readily known by an application without this function has been rejected. Its whole
5195 purpose is to query the typed memory object for attributes that are not user-specified,
5196 but determined by the implementation.
- 5197 4. There should be no separate *mem_alloc()* or *mem_free()* functions. Instead, using
5198 *mmap()* on a typed memory object opened with an “allocate on *mmap()*” flag should be
5199 used to force allocation. These are precisely the semantics defined in the current draft.

5200 • Rationale for no Typed Memory Access Management

5201 The working group had originally defined an additional interface (and an additional kind of
5202 object: typed memory master) to establish and dissolve mappings to typed memory on
5203 behalf of devices or processors which were independent of the operating system and had no
5204 inherent capability to directly establish mappings on their own. This was to have provided
5205 functionality similar to device driver interfaces such as *physio()* and their underlying bus-
5206 specific interfaces (for example, *mballoc()*) which serve to set up and break down DMA
5207 pathways, and derive mapped addresses for use by hardware devices and processor cards.

5208 The ballot group felt that this was beyond the scope of POSIX.1 and its amendments.
5209 Furthermore, the removal of interrupt handling interfaces from a preceding amendment (the
5210 IEEE Std 1003.1d-1999) during its balloting process renders these typed memory access
5211 management interfaces an incomplete solution to portable device management from a user
5212 process; it would be possible to initiate a device transfer to/from typed memory, but
5213 impossible to handle the transfer-complete interrupt in a portable way.

5214 To achieve ballot group consensus, all references to typed memory access management
5215 capabilities were removed. The concept of portable interfaces from a device driver to both
5216 operating system and hardware is being addressed by the Uniform Driver Interface (UDI)
5217 industry forum, with formal standardization deferred until proof of concept and industry-
5218 wide acceptance and implementation.

5219 B.2.8.4 Process Scheduling

5220 IEEE PASC Interpretation 1003.1 #96 has been applied, adding the *pthread_setschedprio()*
5221 function. This was added since previously there was no way for a thread to lower its own
5222 priority without going to the tail of the threads list for its new priority. This capability is
5223 necessary to bound the duration of priority inversion encountered by a thread.

5224 The following portion of the rationale presents models, requirements, and standardization issues
5225 relevant to process scheduling; see also Section B.2.9.4 (on page 168).

5226 In an operating system supporting multiple concurrent processes, the system determines the
5227 order in which processes execute to meet implementation-defined goals. For time-sharing
5228 systems, the goal is to enhance system throughput and promote fairness; the application is
5229 provided with little or no control over this sequencing function. While this is acceptable and
5230 desirable behavior in a time-sharing system, it is inappropriate in a realtime system; realtime
5231 applications must specifically control the execution sequence of their concurrent processes in
5232 order to meet externally defined response requirements.

5233 In IEEE Std 1003.1-2001, the control over process sequencing is provided using a concept of
5234 scheduling policies. These policies, described in detail in this section, define the behavior of the
5235 system whenever processor resources are to be allocated to competing processes. Only the
5236 behavior of the policy is defined; conforming implementations are free to use any mechanism
5237 desired to achieve the described behavior.

5238 • Models

5239 In an operating system supporting multiple concurrent processes, the system determines the
5240 order in which processes execute and might force long-running processes to yield to other
5241 processes at certain intervals. Typically, the scheduling code is executed whenever an event
5242 occurs that might alter the process to be executed next.

5243 The simplest scheduling strategy is a “first-in, first-out” (FIFO) dispatcher. Whenever a
5244 process becomes runnable, it is placed on the end of a ready list. The process at the front of
5245 the ready list is executed until it exits or becomes blocked, at which point it is removed from
5246 the list. This scheduling technique is also known as “run-to-completion” or “run-to-block”.

5247 A natural extension to this scheduling technique is the assignment of a “non-migrating
5248 priority” to each process. This policy differs from strict FIFO scheduling in only one respect:
5249 whenever a process becomes runnable, it is placed at the end of the list of processes runnable
5250 at that priority level. When selecting a process to run, the system always selects the first
5251 process from the highest priority queue with a runnable process. Thus, when a process
5252 becomes unblocked, it will preempt a running process of lower priority without otherwise
5253 altering the ready list. Further, if a process elects to alter its priority, it is removed from the
5254 ready list and reinserted, using its new priority, according to the policy above.

5255 While the above policy might be considered unfriendly in a time-sharing environment in
5256 which multiple users require more balanced resource allocation, it could be ideal in a
5257 realtime environment for several reasons. The most important of these is that it is
5258 deterministic: the highest-priority process is always run and, among processes of equal
5259 priority, the process that has been runnable for the longest time is executed first. Because of
5260 this determinism, cooperating processes can implement more complex scheduling simply by

- 5261 altering their priority. For instance, if processes at a single priority were to reschedule
5262 themselves at fixed time intervals, a time-slice policy would result.
- 5263 In a dedicated operating system in which all processes are well-behaved realtime
5264 applications, non-migrating priority scheduling is sufficient. However, many existing
5265 implementations provide for more complex scheduling policies.
- 5266 IEEE Std 1003.1-2001 specifies a linear scheduling model. In this model, every process in the
5267 system has a priority. The system scheduler always dispatches a process that has the highest
5268 (generally the most time-critical) priority among all runnable processes in the system. As
5269 long as there is only one such process, the dispatching policy is trivial. When multiple
5270 processes of equal priority are eligible to run, they are ordered according to a strict run-to-
5271 completion (FIFO) policy.
- 5272 The priority is represented as a positive integer and is inherited from the parent process. For
5273 processes running under a fixed priority scheduling policy, the priority is never altered
5274 except by an explicit function call.
- 5275 It was determined arbitrarily that larger integers correspond to “higher priorities”.
- 5276 Certain implementations might impose restrictions on the priority ranges to which processes
5277 can be assigned. There also can be restrictions on the set of policies to which processes can be
5278 set.
- 5279 • Requirements
- 5280 Realtime processes require that scheduling be fast and deterministic, and that it guarantees
5281 to preempt lower priority processes.
- 5282 Thus, given the linear scheduling model, realtime processes require that they be run at a
5283 priority that is higher than other processes. Within this framework, realtime processes are
5284 free to yield execution resources to each other in a completely portable and implementation-
5285 defined manner.
- 5286 As there is a generally perceived requirement for processes at the same priority level to share
5287 processor resources more equitably, provisions are made by providing a scheduling policy
5288 (that is, SCHED_RR) intended to provide a timeslice-like facility.
- 5289 **Note:** The following topics assume that low numeric priority implies low scheduling criticality
5290 and *vice versa*.
- 5291 • Rationale for New Interface
- 5292 Realtime applications need to be able to determine when processes will run in relation to
5293 each other. It must be possible to guarantee that a critical process will run whenever it is
5294 runnable; that is, whenever it wants to for as long as it needs. SCHED_FIFO satisfies this
5295 requirement. Additionally, SCHED_RR was defined to meet a realtime requirement for a
5296 well-defined time-sharing policy for processes at the same priority.
- 5297 It would be possible to use the BSD *setpriority()* and *getpriority()* functions by redefining the
5298 meaning of the “nice” parameter according to the scheduling policy currently in use by the
5299 process. The System V *nice()* interface was felt to be undesirable for realtime because it
5300 specifies an adjustment to the “nice” value, rather than setting it to an explicit value.
5301 Realtime applications will usually want to set priority to an explicit value. Also, System V
5302 *nice()* does not allow for changing the priority of another process.
- 5303 With the POSIX.1b interfaces, the traditional “nice” value does not affect the SCHED_FIFO
5304 or SCHED_RR scheduling policies. If a “nice” value is supported, it is implementation-
5305 defined whether it affects the SCHED_OTHER policy.

5306 An important aspect of IEEE Std 1003.1-2001 is the explicit description of the queuing and
5307 preemption rules. It is critical, to achieve deterministic scheduling, that such rules be stated
5308 clearly in IEEE Std 1003.1-2001.

5309 IEEE Std 1003.1-2001 does not address the interaction between priority and swapping. The
5310 issues involved with swapping and virtual memory paging are extremely implementation-
5311 defined and would be nearly impossible to standardize at this point. The proposed
5312 scheduling paradigm, however, fully describes the scheduling behavior of runnable
5313 processes, of which one criterion is that the working set be resident in memory. Assuming
5314 the existence of a portable interface for locking portions of a process in memory, paging
5315 behavior need not affect the scheduling of realtime processes.

5316 IEEE Std 1003.1-2001 also does not address the priorities of “system” processes. In general,
5317 these processes should always execute in low-priority ranges to avoid conflict with other
5318 realtime processes. Implementations should document the priority ranges in which system
5319 processes run.

5320 The default scheduling policy is not defined. The effect of I/O interrupts and other system
5321 processing activities is not defined. The temporary lending of priority from one process to
5322 another (such as for the purposes of affecting freeing resources) by the system is not
5323 addressed. Preemption of resources is not addressed. Restrictions on the ability of a process
5324 to affect other processes beyond a certain level (influence levels) is not addressed.

5325 The rationale used to justify the simple time-quantum scheduler is that it is common practice
5326 to depend upon this type of scheduling to ensure “fair” distribution of processor resources
5327 among portions of the application that must interoperate in a serial fashion. Note that
5328 IEEE Std 1003.1-2001 is silent with respect to the setting of this time quantum, or whether it is
5329 a system-wide value or a per-process value, although it appears that the prevailing realtime
5330 practice is for it to be a system-wide value.

5331 In a system with N processes at a given priority, all processor-bound, in which the time
5332 quantum is equal for all processes at a specific priority level, the following assumptions are
5333 made of such a scheduling policy:

- 5334 1. A time quantum Q exists and the current process will own control of the processor for
5335 at least a duration of Q and will have the processor for a duration of Q .
- 5336 2. The N th process at that priority will control a processor within a duration of $(N-1) \times Q$.

5337 These assumptions are necessary to provide equal access to the processor and bounded
5338 response from the application.

5339 The assumptions hold for the described scheduling policy only if no system overhead, such
5340 as interrupt servicing, is present. If the interrupt servicing load is non-zero, then one of the
5341 two assumptions becomes fallacious, based upon how Q is measured by the system.

5342 If Q is measured by clock time, then the assumption that the process obtains a duration Q
5343 processor time is false if interrupt overhead exists. Indeed, a scenario can be constructed with
5344 N processes in which a single process undergoes complete processor starvation if a
5345 peripheral device, such as an analog-to-digital converter, generates significant interrupt
5346 activity periodically with a period of $N \times Q$.

5347 If Q is measured as actual processor time, then the assumption that the N th process runs in
5348 within the duration $(N-1) \times Q$ is false.

5349 It should be noted that SCHED_FIFO suffers from interrupt-based delay as well. However,
5350 for SCHED_FIFO, the implied response of the system is “as soon as possible”, so that the
5351 interrupt load for this case is a vendor selection and not a compliance issue.

5352 With this in mind, it is necessary either to complete the definition by including bounds on the
5353 interrupt load, or to modify the assumptions that can be made about the scheduling policy.

5354 Since the motivation of inclusion of the policy is common usage, and since current
5355 applications do not enjoy the luxury of bounded interrupt load, item (2) above is sufficient to
5356 express existing application needs and is less restrictive in the standard definition. No
5357 difference in interface is necessary.

5358 In an implementation in which the time quantum is equal for all processes at a specific
5359 priority, our assumptions can then be restated as:

5360 — A time quantum Q exists, and a processor-bound process will be rescheduled after a
5361 duration of, at most, Q . Time quantum Q may be defined in either wall clock time or
5362 execution time.

5363 — In general, the N th process of a priority level should wait no longer than $(N-1) \times Q$ time
5364 to execute, assuming no processes exist at higher priority levels.

5365 — No process should wait indefinitely.

5366 For implementations supporting per-process time quanta, these assumptions can be readily
5367 extended.

5368 **Sporadic Server Scheduling Policy**

5369 The sporadic server is a mechanism defined for scheduling aperiodic activities in time-critical
5370 realtime systems. This mechanism reserves a certain bounded amount of execution capacity for
5371 processing aperiodic events at a high priority level. Any aperiodic events that cannot be
5372 processed within the bounded amount of execution capacity are executed in the background at a
5373 low priority level. Thus, a certain amount of execution capacity can be guaranteed to be
5374 available for processing periodic tasks, even under burst conditions in the arrival of aperiodic
5375 processing requests (that is, a large number of requests in a short time interval). The sporadic
5376 server also simplifies the schedulability analysis of the realtime system, because it allows
5377 aperiodic processes or threads to be treated as if they were periodic. The sporadic server was
5378 first described by Sprunt, et al.

5379 The key concept of the sporadic server is to provide and limit a certain amount of computation
5380 capacity for processing aperiodic events at their assigned normal priority, during a time interval
5381 called the “replenishment period”. Once the entity controlled by the sporadic server mechanism
5382 is initialized with its period and execution-time budget attributes, it preserves its execution
5383 capacity until an aperiodic request arrives. The request will be serviced (if there are no higher
5384 priority activities pending) as long as there is execution capacity left. If the request is completed,
5385 the actual execution time used to service it is subtracted from the capacity, and a replenishment
5386 of this amount of execution time is scheduled to happen one replenishment period after the
5387 arrival of the aperiodic request. If the request is not completed, because there is no execution
5388 capacity left, then the aperiodic process or thread is assigned a lower background priority. For
5389 each portion of consumed execution capacity the execution time used is replenished after one
5390 replenishment period. At the time of replenishment, if the sporadic server was executing at a
5391 background priority level, its priority is elevated to the normal level. Other similar
5392 replenishment policies have been defined, but the one presented here represents a compromise
5393 between efficiency and implementation complexity.

5394 The interface that appears in this section defines a new scheduling policy for threads and
5395 processes that behaves according to the rules of the sporadic server mechanism. Scheduling
5396 attributes are defined and functions are provided to allow the user to set and get the parameters
5397 that control the scheduling behavior of this mechanism, namely the normal and low priority, the
5398 replenishment period, the maximum number of pending replenishment operations, and the

- 5399 initial execution-time budget.
- 5400 • Scheduling Aperiodic Activities
- 5401 Virtually all realtime applications are required to process aperiodic activities. In many cases,
5402 there are tight timing constraints that the response to the aperiodic events must meet. Usual
5403 timing requirements imposed on the response to these events are:
- 5404 — The effects of an aperiodic activity on the response time of lower priority activities must
5405 be controllable and predictable.
 - 5406 — The system must provide the fastest possible response time to aperiodic events.
 - 5407 — It must be possible to take advantage of all the available processing bandwidth not
5408 needed by time-critical activities to enhance average-case response times to aperiodic
5409 events.
- 5410 Traditional methods for scheduling aperiodic activities are background processing, polling
5411 tasks, and direct event execution:
- 5412 — Background processing consists of assigning a very low priority to the processing of
5413 aperiodic events. It utilizes all the available bandwidth in the system that has not been
5414 consumed by higher priority threads. However, it is very difficult, or impossible, to meet
5415 requirements on average-case response time, because the aperiodic entity has to wait for
5416 the execution of all other entities which have higher priority.
 - 5417 — Polling consists of creating a periodic process or thread for servicing aperiodic requests.
5418 At regular intervals, the polling entity is started and its services accumulated pending
5419 aperiodic requests. If no aperiodic requests are pending, the polling entity suspends itself
5420 until its next period. Polling allows the aperiodic requests to be processed at a higher
5421 priority level. However, worst and average-case response times of polling entities are a
5422 direct function of the polling period, and there is execution overhead for each polling
5423 period, even if no event has arrived. If the deadline of the aperiodic activity is short
5424 compared to the inter-arrival time, the polling frequency must be increased to guarantee
5425 meeting the deadline. For this case, the increase in frequency can dramatically reduce the
5426 efficiency of the system and, therefore, its capacity to meet all deadlines. Yet, polling
5427 represents a good way to handle a large class of practical problems because it preserves
5428 system predictability, and because the amortized overhead drops as load increases.
 - 5429 — Direct event execution consists of executing the aperiodic events at a high fixed-priority
5430 level. Typically, the aperiodic event is processed by an interrupt service routine as soon as
5431 it arrives. This technique provides predictable response times for aperiodic events, but
5432 makes the response times of all lower priority activities completely unpredictable under
5433 burst arrival conditions. Therefore, if the density of aperiodic event arrivals is
5434 unbounded, it may be a dangerous technique for time-critical systems. Yet, for those cases
5435 in which the physics of the system imposes a bound on the event arrival rate, it is
5436 probably the most efficient technique.
 - 5437 — The sporadic server scheduling algorithm combines the predictability of the polling
5438 approach with the short response times of the direct event execution. Thus, it allows
5439 systems to meet an important class of application requirements that cannot be met by
5440 using the traditional approaches. Multiple sporadic servers with different attributes can
5441 be applied to the scheduling of multiple classes of aperiodic events, each with different
5442 kinds of timing requirements, such as individual deadlines, average response times, and
5443 so on. It also has many other interesting applications for realtime, such as scheduling
5444 producer/consumer tasks in time-critical systems, limiting the effects of faults on the
5445 estimation of task execution-time requirements, and so on.

- 5446
- Existing Practice
- 5447 The sporadic server has been used in different kinds of applications, including military
5448 avionics, robot control systems, industrial automation systems, and so on. There are
5449 examples of many systems that cannot be successfully scheduled using the classic
5450 approaches, such as direct event execution, or polling, and are schedulable using a sporadic
5451 server scheduler. The sporadic server algorithm itself can successfully schedule all systems
5452 scheduled with direct event execution or polling.
- 5453 The sporadic server scheduling policy has been implemented as a commercial product in the
5454 run-time system of the Verdex Ada compiler. There are also many applications that have
5455 used a much less efficient application-level sporadic server. These realtime applications
5456 would benefit from a sporadic server scheduler implemented at the scheduler level.
- Library-Level *versus* Kernel-Level Implementation
- 5458 The sporadic server interface described in this section requires the sporadic server policy to
5459 be implemented at the same level as the scheduler. This means that the process sporadic
5460 server must be implemented at the kernel level and the thread sporadic server policy
5461 implemented at the same level as the thread scheduler; that is, kernel or library level.
- 5462 In an earlier interface for the sporadic server, this mechanism was implementable at a
5463 different level than the scheduler. This feature allowed the implementor to choose between
5464 an efficient scheduler-level implementation, or a simpler user or library-level
5465 implementation. However, the working group considered that this interface made the use of
5466 sporadic servers more complex, and that library-level implementations would lack some of
5467 the important functionality of the sporadic server, namely the limitation of the actual
5468 execution time of aperiodic activities. The working group also felt that the interface
5469 described in this chapter does not preclude library-level implementations of threads intended
5470 to provide efficient low-overhead scheduling for those threads that are not scheduled under
5471 the sporadic server policy.
- Range of Scheduling Priorities
- 5473 Each of the scheduling policies supported in IEEE Std 1003.1-2001 has an associated range of
5474 priorities. The priority ranges for each policy might or might not overlap with the priority
5475 ranges of other policies. For time-critical realtime applications it is usual for periodic and
5476 aperiodic activities to be scheduled together in the same processor. Periodic activities will
5477 usually be scheduled using the SCHED_FIFO scheduling policy, while aperiodic activities
5478 may be scheduled using SCHED_SPORADIC. Since the application developer will require
5479 complete control over the relative priorities of these activities in order to meet his timing
5480 requirements, it would be desirable for the priority ranges of SCHED_FIFO and
5481 SCHED_SPORADIC to overlap completely. Therefore, although IEEE Std 1003.1-2001 does
5482 not require any particular relationship between the different priority ranges, it is
5483 recommended that these two ranges should coincide.
- Dynamically Setting the Sporadic Server Policy
- 5485 Several members of the working group requested that implementations should not be
5486 required to support dynamically setting the sporadic server scheduling policy for a thread.
5487 The reason is that this policy may have a high overhead for library-level implementations of
5488 threads, and if threads are allowed to dynamically set this policy, this overhead can be
5489 experienced even if the thread does not use that policy. By disallowing the dynamic setting
5490 of the sporadic server scheduling policy, these implementations can accomplish efficient
5491 scheduling for threads using other policies. If a strictly conforming application needs to use
5492 the sporadic server policy, and is therefore willing to pay the overhead, it must set this policy
5493 at the time of thread creation.

5494 • Limitation of the Number of Pending Replenishments

5495 The number of simultaneously pending replenishment operations must be limited for each
 5496 sporadic server for two reasons: an unlimited number of replenishment operations would
 5497 need an unlimited number of system resources to store all the pending replenishment
 5498 operations; on the other hand, in some implementations each replenishment operation will
 5499 represent a source of priority inversion (just for the duration of the replenishment operation)
 5500 and thus, the maximum amount of replenishments must be bounded to guarantee bounded
 5501 response times. The way in which the number of replenishments is bounded is by lowering
 5502 the priority of the sporadic server to *sched_ss_low_priority* when the number of pending
 5503 replenishments has reached its limit. In this way, no new replenishments are scheduled until
 5504 the number of pending replenishments decreases.

5505 In the sporadic server scheduling policy defined in IEEE Std 1003.1-2001, the application can
 5506 specify the maximum number of pending replenishment operations for a single sporadic
 5507 server, by setting the value of the *sched_ss_max_repl* scheduling parameter. This value must
 5508 be between one and {SS_REPL_MAX}, which is a maximum limit imposed by the
 5509 implementation. The limit {SS_REPL_MAX} must be greater than or equal to
 5510 {_POSIX_SS_REPL_MAX}, which is defined to be four in IEEE Std 1003.1-2001. The minimum
 5511 limit of four was chosen so that an application can at least guarantee that four different
 5512 aperiodic events can be processed during each interval of length equal to the replenishment
 5513 period.

5514 *B.2.8.5 Clocks and Timers*

5515 • Clocks

5516 IEEE Std 1003.1-2001 and the ISO C standard both define functions for obtaining system time.
 5517 Implicit behind these functions is a mechanism for measuring passage of time. This
 5518 specification makes this mechanism explicit and calls it a clock. The *CLOCK_REALTIME*
 5519 clock required by IEEE Std 1003.1-2001 is a higher resolution version of the clock that
 5520 maintains POSIX.1 system time. This is a “system-wide” clock, in that it is visible to all
 5521 processes and, were it possible for multiple processes to all read the clock at the same time,
 5522 they would see the same value.

5523 An extensible interface was defined, with the ability for implementations to define additional
 5524 clocks. This was done because of the observation that many realtime platforms support
 5525 multiple clocks, and it was desired to fit this model within the standard interface. But
 5526 implementation-defined clocks need not represent actual hardware devices, nor are they
 5527 necessarily system-wide.

5528 • Timers

5529 Two timer types are required for a system to support realtime applications:

5530 1. One-shot

5531 A one-shot timer is a timer that is armed with an initial expiration time, either relative
 5532 to the current time or at an absolute time (based on some timing base, such as time in
 5533 seconds and nanoseconds since the Epoch). The timer expires once and then is
 5534 disarmed. With the specified facilities, this is accomplished by setting the *it_value*
 5535 member of the *value* argument to the desired expiration time and the *it_interval* member
 5536 to zero.

5537 2. Periodic

5538 A periodic timer is a timer that is armed with an initial expiration time, again either
 5539 relative or absolute, and a repetition interval. When the initial expiration occurs, the

5540 timer is reloaded with the repetition interval and continues counting. With the
 5541 specified facilities, this is accomplished by setting the *it_value* member of the *value*
 5542 argument to the desired initial expiration time and the *it_interval* member to the desired
 5543 repetition interval.

5544 For both of these types of timers, the time of the initial timer expiration can be specified in
 5545 two ways:

- 5546 1. Relative (to the current time)
- 5547 2. Absolute

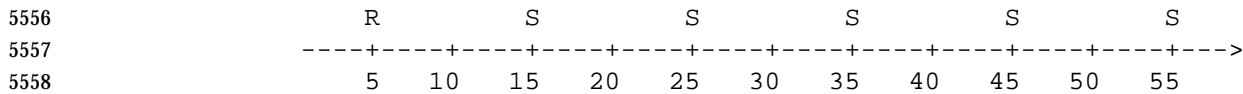
5548 • Examples of Using Realtime Timers

5549 In the diagrams below, *S* indicates a program schedule, *R* shows a schedule method request,
 5550 and *E* suggests an internal operating system event.

5551 — Periodic Timer: Data Logging

5552 During an experiment, it might be necessary to log realtime data periodically to an
 5553 internal buffer or to a mass storage device. With a periodic scheduling method, a logging
 5554 module can be started automatically at fixed time intervals to log the data.

5555 Program schedule is requested every 10 seconds.



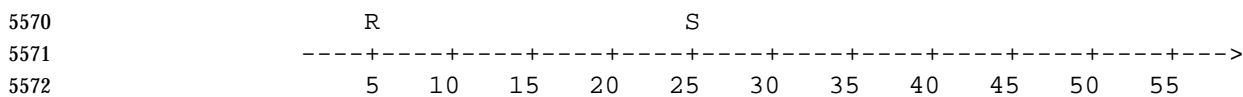
5559 [Time (in Seconds)]

5560 To achieve this type of scheduling using the specified facilities, one would allocate a per-
 5561 process timer based on clock ID `CLOCK_REALTIME`. Then the timer would be armed via
 5562 a call to `timer_settime()` with the `TIMER_ABSTIME` flag reset, and with an initial
 5563 expiration value and a repetition interval of 10 seconds.

5564 — One-shot Timer (Relative Time): Device Initialization

5565 In an emission test environment, large sample bags are used to capture the exhaust from
 5566 a vehicle. The exhaust is purged from these bags before each and every test. With a one-
 5567 shot timer, a module could initiate the purge function and then suspend itself for a
 5568 predetermined period of time while the sample bags are prepared.

5569 Program schedule requested 20 seconds after call is issued.



5573 [Time (in Seconds)]

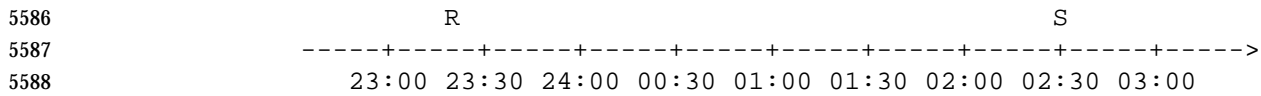
5574 To achieve this type of scheduling using the specified facilities, one would allocate a per-
 5575 process timer based on clock ID `CLOCK_REALTIME`. Then the timer would be armed via
 5576 a call to `timer_settime()` with the `TIMER_ABSTIME` flag reset, and with an initial
 5577 expiration value of 20 seconds and a repetition interval of zero.

5578 Note that if the program wishes merely to suspend itself for the specified interval, it
 5579 could more easily use `nanosleep()`.

5580 — One-shot Timer (Absolute Time): Data Transmission

5581 The results from an experiment are often moved to a different system within a network
 5582 for postprocessing or archiving. With an absolute one-shot timer, a module that moves
 5583 data from a test-cell computer to a host computer can be automatically scheduled on a
 5584 daily basis.

5585 Program schedule requested for 2:30 a.m.



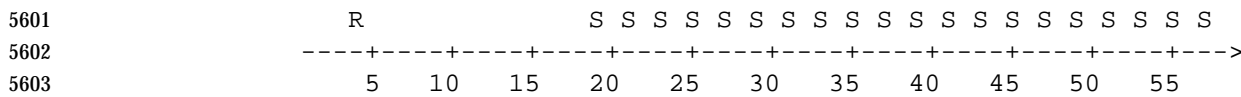
5589 [Time of Day]

5590 To achieve this type of scheduling using the specified facilities, a per-process timer would
 5591 be allocated based on clock ID CLOCK_REALTIME. Then the timer would be armed via
 5592 a call to *timer_settime()* with the TIMER_ABSTIME flag set, and an initial expiration value
 5593 equal to 2:30 a.m. of the next day.

5594 — Periodic Timer (Relative Time): Signal Stabilization

5595 Some measurement devices, such as emission analyzers, do not respond instantaneously
 5596 to an introduced sample. With a periodic timer with a relative initial expiration time, a
 5597 module that introduces a sample and records the average response could suspend itself
 5598 for a predetermined period of time while the signal is stabilized and then sample at a
 5599 fixed rate.

5600 Program schedule requested 15 seconds after call is issued and every 2 seconds thereafter.



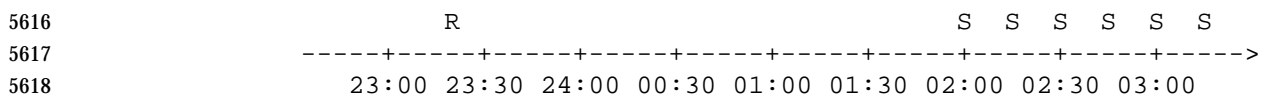
5604 [Time (in Seconds)]

5605 To achieve this type of scheduling using the specified facilities, one would allocate a per-
 5606 process timer based on clock ID CLOCK_REALTIME. Then the timer would be armed via
 5607 a call to *timer_settime()* with TIMER_ABSTIME flag reset, and with an initial expiration
 5608 value of 15 seconds and a repetition interval of 2 seconds.

5609 — Periodic Timer (Absolute Time): Work Shift-related Processing

5610 Resource utilization data is useful when time to perform experiments is being scheduled
 5611 at a facility. With a periodic timer with an absolute initial expiration time, a module can
 5612 be scheduled at the beginning of a work shift to gather resource utilization data
 5613 throughout the shift. This data can be used to allocate resources effectively to minimize
 5614 bottlenecks and delays and maximize facility throughput.

5615 Program schedule requested for 2:00 a.m. and every 15 minutes thereafter.



5619 [Time of Day]

5620 To achieve this type of scheduling using the specified facilities, one would allocate a per-
 5621 process timer based on clock ID CLOCK_REALTIME. Then the timer would be armed via
 5622 a call to *timer_settime()* with TIMER_ABSTIME flag set, and with an initial expiration
 5623 value equal to 2:00 a.m. and a repetition interval equal to 15 minutes.

5624 • Relationship of Timers to Clocks

5625 The relationship between clocks and timers armed with an absolute time is straightforward:
 5626 a timer expiration signal is requested when the associated clock reaches or exceeds the
 5627 specified time. The relationship between clocks and timers armed with a relative time (an
 5628 interval) is less obvious, but not unintuitive. In this case, a timer expiration signal is
 5629 requested when the specified interval, *as measured by the associated clock*, has passed. For the
 5630 required CLOCK_REALTIME clock, this allows timer expiration signals to be requested at
 5631 specified “wall clock” times (absolute), or when a specified interval of “realtime” has passed
 5632 (relative). For an implementation-defined clock—say, a process virtual time clock—timer
 5633 expirations could be requested when the process has used a specified total amount of virtual
 5634 time (absolute), or when it has used a specified *additional* amount of virtual time (relative).

5635 The interfaces also allow flexibility in the implementation of the functions. For example, an
 5636 implementation could convert all absolute times to intervals by subtracting the clock value at
 5637 the time of the call from the requested expiration time and “counting down” at the
 5638 supported resolution. Or it could convert all relative times to absolute expiration time by
 5639 adding in the clock value at the time of the call and comparing the clock value to the
 5640 expiration time at the supported resolution. Or it might even choose to maintain absolute
 5641 times as absolute and compare them to the clock value at the supported resolution for
 5642 absolute timers, and maintain relative times as intervals and count them down at the
 5643 resolution supported for relative timers. The choice will be driven by efficiency
 5644 considerations and the underlying hardware or software clock implementation.

5645 • Data Definitions for Clocks and Timers

5646 IEEE Std 1003.1-2001 uses a time representation capable of supporting nanosecond resolution
 5647 timers for the following reasons:

- 5648 — To enable IEEE Std 1003.1-2001 to represent those computer systems already using
 5649 nanosecond or submicrosecond resolution clocks.
- 5650 — To accommodate those per-process timers that might need nanoseconds to specify an
 5651 absolute value of system-wide clocks, even though the resolution of the per-process timer
 5652 may only be milliseconds, or *vice versa*.
- 5653 — Because the number of nanoseconds in a second can be represented in 32 bits.

5654 Time values are represented in the **timespec** structure. The *tv_sec* member is of type **time_t**
 5655 so that this member is compatible with time values used by POSIX.1 functions and the ISO C
 5656 standard. The *tv_nsec* member is a **signed long** in order to simplify and clarify code that
 5657 decrements or finds differences of time values. Note that because 1 billion (number of
 5658 nanoseconds per second) is less than half of the value representable by a signed 32-bit value,
 5659 it is always possible to add two valid fractional seconds represented as integral nanoseconds
 5660 without overflowing the signed 32-bit value.

5661 A maximum allowable resolution for the CLOCK_REALTIME clock of 20 ms (1/50 seconds)
 5662 was chosen to allow line frequency clocks in European countries to be conforming. 60 Hz
 5663 clocks in the U.S. will also be conforming, as will finer granularity clocks, although a Strictly
 5664 Conforming Application cannot assume a granularity of less than 20 ms (1/50 seconds).

5665 The minimum allowable maximum time allowed for the CLOCK_REALTIME clock and the
 5666 function *nanosleep()*, and timers created with *clock_id*=CLOCK_REALTIME, is determined by
 5667 the fact that the *tv_sec* member is of type **time_t**.

5668 IEEE Std 1003.1-2001 specifies that timer expirations must not be delivered early, and
 5669 *nanosleep()* must not return early due to quantization error. IEEE Std 1003.1-2001 discusses
 5670 the various implementations of *alarm()* in the rationale and states that implementations that

5671 do not allow alarm signals to occur early are the most appropriate, but refrained from
5672 mandating this behavior. Because of the importance of predictability to realtime applications,
5673 IEEE Std 1003.1-2001 takes a stronger stance.

5674 The developers of IEEE Std 1003.1-2001 considered using a time representation that differs
5675 from POSIX.1b in the second 32 bit of the 64-bit value. Whereas POSIX.1b defines this field
5676 as a fractional second in nanoseconds, the other methodology defines this as a binary fraction
5677 of one second, with the radix point assumed before the most significant bit.

5678 POSIX.1b is a software, source-level standard and most of the benefits of the alternate
5679 representation are enjoyed by hardware implementations of clocks and algorithms. It was
5680 felt that mandating this format for POSIX.1b clocks and timers would unnecessarily burden
5681 the application writer with writing, possibly non-portable, multiple precision arithmetic
5682 packages to perform conversion between binary fractions and integral units such as
5683 nanoseconds, milliseconds, and so on.

5684 **Rationale for the Monotonic Clock**

5685 For those applications that use time services to achieve realtime behavior, changing the value of
5686 the clock on which these services rely may cause erroneous timing behavior. For these
5687 applications, it is necessary to have a monotonic clock which cannot run backwards, and which
5688 has a maximum clock jump that is required to be documented by the implementation.
5689 Additionally, it is desirable (but not required by IEEE Std 1003.1-2001) that the monotonic clock
5690 increases its value uniformly. This clock should not be affected by changes to the system time;
5691 for example, to synchronize the clock with an external source or to account for leap seconds.
5692 Such changes would cause errors in the measurement of time intervals for those time services
5693 that use the absolute value of the clock.

5694 One could argue that by defining the behavior of time services when the value of a clock is
5695 changed, deterministic realtime behavior can be achieved. For example, one could specify that
5696 relative time services should be unaffected by changes in the value of a clock. However, there
5697 are time services that are based upon an absolute time, but that are essentially intended as
5698 relative time services. For example, *pthread_cond_timedwait()* uses an absolute time to allow it to
5699 wake up after the required interval despite spurious wakeups. Although sometimes the
5700 *pthread_cond_timedwait()* timeouts are absolute in nature, there are many occasions in which
5701 they are relative, and their absolute value is determined from the current time plus a relative
5702 time interval. In this latter case, if the clock changes while the thread is waiting, the wait interval
5703 will not be the expected length. If a *pthread_cond_timedwait()* function were created that would
5704 take a relative time, it would not solve the problem because to retain the intended “deadline” a
5705 thread would need to compensate for latency due to the spurious wakeup, and preemption
5706 between wakeup and the next wait.

5707 The solution is to create a new monotonic clock, whose value does not change except for the
5708 regular ticking of the clock, and use this clock for implementing the various relative timeouts
5709 that appear in the different POSIX interfaces, as well as allow *pthread_cond_timedwait()* to choose
5710 this new clock for its timeout. A new *clock_nanosleep()* function is created to allow an application
5711 to take advantage of this newly defined clock. Notice that the monotonic clock may be
5712 implemented using the same hardware clock as the system clock.

5713 Relative timeouts for *sigtimedwait()* and *aio_suspend()* have been redefined to use the monotonic
5714 clock, if present. The *alarm()* function has not been redefined, because the same effect but with
5715 better resolution can be achieved by creating a timer (for which the appropriate clock may be
5716 chosen).

5717 The *pthread_cond_timedwait()* function has been treated in a different way, compared to other
5718 functions with absolute timeouts, because it is used to wait for an event, and thus it may have a

5719 deadline, while the other timeouts are generally used as an error recovery mechanism, and for
5720 them the use of the monotonic clock is not so important. Since the desired timeout for the
5721 *pthread_cond_timedwait()* function may either be a relative interval or an absolute time of day
5722 deadline, a new initialization attribute has been created for condition variables to specify the
5723 clock that is used for measuring the timeout in a call to *pthread_cond_timedwait()*. In this way, if
5724 a relative timeout is desired, the monotonic clock will be used; if an absolute deadline is required
5725 instead, the `CLOCK_REALTIME` or another appropriate clock may be used. This capability has
5726 not been added to other functions with absolute timeouts because for those functions the
5727 expected use of the timeout is mostly to prevent errors, and not so often to meet precise
5728 deadlines. As a consequence, the complexity of adding this capability is not justified by its
5729 perceived application usage.

5730 The *nanosleep()* function has not been modified with the introduction of the monotonic clock.
5731 Instead, a new *clock_nanosleep()* function has been created, in which the desired clock may be
5732 specified in the function call.

5733 • History of Resolution Issues

5734 Due to the shift from relative to absolute timeouts in IEEE Std 1003.1d-1999, the amendments
5735 to the *sem_timedwait()*, *pthread_mutex_timedlock()*, *mq_timedreceive()*, and *mq_timedsend()*
5736 functions of that standard have been removed. Those amendments specified that
5737 `CLOCK_MONOTONIC` would be used for the (relative) timeouts if the Monotonic Clock
5738 option was supported.

5739 Having these functions continue to be tied solely to `CLOCK_MONOTONIC` would not
5740 work. Since the absolute value of a time value obtained from `CLOCK_MONOTONIC` is
5741 unspecified, under the absolute timeouts interface, applications would behave differently
5742 depending on whether the Monotonic Clock option was supported or not (because the
5743 absolute value of the clock would have different meanings in either case).

5744 Two options were considered:

- 5745 1. Leave the current behavior unchanged, which specifies the `CLOCK_REALTIME` clock
5746 for these (absolute) timeouts, to allow portability of applications between
5747 implementations supporting or not the Monotonic Clock option.
- 5748 2. Modify these functions in the way that *pthread_cond_timedwait()* was modified to allow
5749 a choice of clock, so that an application could use `CLOCK_REALTIME` when it is trying
5750 to achieve an absolute timeout and `CLOCK_MONOTONIC` when it is trying to achieve
5751 a relative timeout.

5752 It was decided that the features of `CLOCK_MONOTONIC` are not as critical to these
5753 functions as they are to *pthread_cond_timedwait()*. The *pthread_cond_timedwait()* function is
5754 given a relative timeout; the timeout may represent a deadline for an event. When these
5755 functions are given relative timeouts, the timeouts are typically for error recovery purposes
5756 and need not be so precise.

5757 Therefore, it was decided that these functions should be tied to `CLOCK_REALTIME` and not
5758 complicated by being given a choice of clock.

5759 **Execution Time Monitoring**

5760 • Introduction

5761 The main goals of the execution time monitoring facilities defined in this chapter are to
5762 measure the execution time of processes and threads and to allow an application to establish
5763 CPU time limits for these entities.

5764 The analysis phase of time-critical realtime systems often relies on the measurement of
5765 execution times of individual threads or processes to determine whether the timing
5766 requirements will be met. Also, performance analysis techniques for soft deadline realtime
5767 systems rely heavily on the determination of these execution times. The execution time
5768 monitoring functions provide application developers with the ability to measure these
5769 execution times online and open the possibility of dynamic execution-time analysis and
5770 system reconfiguration, if required.

5771 The second goal of allowing an application to establish execution time limits for individual
5772 processes or threads and detecting when they overrun allows program robustness to be
5773 increased by enabling online checking of the execution times.

5774 If errors are detected—possibly because of erroneous program constructs, the existence of
5775 errors in the analysis phase, or a burst of event arrivals—online detection and recovery is
5776 possible in a portable way. This feature can be extremely important for many time-critical
5777 applications. Other applications require trapping CPU-time errors as a normal way to exit an
5778 algorithm; for instance, some realtime artificial intelligence applications trigger a number of
5779 independent inference processes of varying accuracy and speed, limit how long they can run,
5780 and pick the best answer available when time runs out. In many periodic systems, overrun
5781 processes are simply restarted in the next resource period, after necessary end-of-period
5782 actions have been taken. This allows algorithms that are inherently data-dependent to be
5783 made predictable.

5784 The interface that appears in this chapter defines a new type of clock, the CPU-time clock,
5785 which measures execution time. Each process or thread can invoke the clock and timer
5786 functions defined in POSIX.1 to use them. Functions are also provided to access the CPU-
5787 time clock of other processes or threads to enable remote monitoring of these clocks.
5788 Monitoring of threads of other processes is not supported, since these threads are not visible
5789 from outside of their own process with the interfaces defined in POSIX.1.

5790 • Execution Time Monitoring Interface

5791 The clock and timer interface defined in POSIX.1 historically only defined one clock, which
5792 measures wall-clock time. The requirements for measuring execution time of processes and
5793 threads, and setting limits to their execution time by detecting when they overrun, can be
5794 accomplished with that interface if a new kind of clock is defined. These new clocks measure
5795 execution time, and one is associated with each process and with each thread. The clock
5796 functions currently defined in POSIX.1 can be used to read and set these CPU-time clocks,
5797 and timers can be created using these clocks as their timing base. These timers can then be
5798 used to send a signal when some specified execution time has been exceeded. The CPU-time
5799 clocks of each process or thread can be accessed by using the symbols
5800 CLOCK_PROCESS_CPUTIME_ID or CLOCK_THREAD_CPUTIME_ID.

5801 The clock and timer interface defined in POSIX.1 and extended with the new kind of CPU-
5802 time clock would only allow processes or threads to access their own CPU-time clocks.
5803 However, many realtime systems require the possibility of monitoring the execution time of
5804 processes or threads from independent monitoring entities. In order to allow applications to
5805 construct independent monitoring entities that do not require cooperation from or
5806 modification of the monitored entities, two functions have been added: *clock_getcpucllockid()*,

5807 for accessing CPU-time clocks of other processes, and *pthread_getcpuclockid()*, for accessing
5808 CPU-time clocks of other threads. These functions return the clock identifier associated with
5809 the process or thread specified in the call. These clock IDs can then be used in the rest of the
5810 clock function calls.

5811 The clocks accessed through these functions could also be used as a timing base for the
5812 creation of timers, thereby allowing independent monitoring entities to limit the CPU time
5813 consumed by other entities. However, this possibility would imply additional complexity
5814 and overhead because of the need to maintain a timer queue for each process or thread, to
5815 store the different expiration times associated with timers created by different processes or
5816 threads. The working group decided this additional overhead was not justified by
5817 application requirements. Therefore, creation of timers attached to the CPU-time clocks of
5818 other processes or threads has been specified as implementation-defined.

5819 • Overhead Considerations

5820 The measurement of execution time may introduce additional overhead in the thread
5821 scheduling, because of the need to keep track of the time consumed by each of these entities.
5822 In library-level implementations of threads, the efficiency of scheduling could be somehow
5823 compromised because of the need to make a kernel call, at each context switch, to read the
5824 process CPU-time clock. Consequently, a thread creation attribute called *cpu-clock-*
5825 *requirement* was defined, to allow threads to disconnect their respective CPU-time clocks.
5826 However, the Ballot Group considered that this attribute itself introduced some overhead,
5827 and that in current implementations it was not worth the effort. Therefore, the attribute was
5828 deleted, and thus thread CPU-time clocks are required for all threads if the Thread CPU-Time
5829 Clocks option is supported.

5830 • Accuracy of CPU-Time Clocks

5831 The mechanism used to measure the execution time of processes and threads is specified in
5832 IEEE Std 1003.1-2001 as implementation-defined. The reason for this is that both the
5833 underlying hardware and the implementation architecture have a very strong influence on
5834 the accuracy achievable for measuring CPU time. For some implementations, the
5835 specification of strict accuracy requirements would represent very large overheads, or even
5836 the impossibility of being implemented.

5837 Since the mechanism for measuring execution time is implementation-defined, realtime
5838 applications will be able to take advantage of accurate implementations using a portable
5839 interface. Of course, strictly conforming applications cannot rely on any particular degree of
5840 accuracy, in the same way as they cannot rely on a very accurate measurement of wall clock
5841 time. There will always exist applications whose accuracy or efficiency requirements on the
5842 implementation are more rigid than the values defined in IEEE Std 1003.1-2001 or any other
5843 standard.

5844 In any case, there is a minimum set of characteristics that realtime applications would expect
5845 from most implementations. One such characteristic is that the sum of all the execution times
5846 of all the threads in a process equals the process execution time, when no CPU-time clocks
5847 are disabled. This need not always be the case because implementations may differ in how
5848 they account for time during context switches. Another characteristic is that the sum of the
5849 execution times of all processes in a system equals the number of processors, multiplied by
5850 the elapsed time, assuming that no processor is idle during that elapsed time. However, in
5851 some implementations it might not be possible to relate CPU time to elapsed time. For
5852 example, in a heterogeneous multi-processor system in which each processor runs at a
5853 different speed, an implementation may choose to define each “second” of CPU time to be a
5854 certain number of “cycles” that a CPU has executed.

- 5855 • Existing Practice
- 5856 Measuring and limiting the execution time of each concurrent activity are common features
5857 of most industrial implementations of realtime systems. Almost all critical realtime systems
5858 are currently built upon a cyclic executive. With this approach, a regular timer interrupt kicks
5859 off the next sequence of computations. It also checks that the current sequence has
5860 completed. If it has not, then some error recovery action can be undertaken (or at least an
5861 overrun is avoided). Current software engineering principles and the increasing complexity
5862 of software are driving application developers to implement these systems on multi-
5863 threaded or multi-process operating systems. Therefore, if a POSIX operating system is to be
5864 used for this type of application, then it must offer the same level of protection.
- 5865 Execution time clocks are also common in most UNIX implementations, although these
5866 clocks usually have requirements different from those of realtime applications. The POSIX.1
5867 *times()* function supports the measurement of the execution time of the calling process, and
5868 its terminated child processes. This execution time is measured in clock ticks and is supplied
5869 as two different values with the user and system execution times, respectively. BSD supports
5870 the function *getrusage()*, which allows the calling process to get information about the
5871 resources used by itself and/or all of its terminated child processes. The resource usage
5872 includes user and system CPU time. Some UNIX systems have options to specify high
5873 resolution (up to one microsecond) CPU-time clocks using the *times()* or the *getrusage()*
5874 functions.
- 5875 The *times()* and *getrusage()* interfaces do not meet important realtime requirements, such as
5876 the possibility of monitoring execution time from a different process or thread, or the
5877 possibility of detecting an execution time overrun. The latter requirement is supported in
5878 some UNIX implementations that are able to send a signal when the execution time of a
5879 process has exceeded some specified value. For example, BSD defines the functions
5880 *getitimer()* and *setitimer()*, which can operate either on a realtime clock (wall-clock), or on
5881 virtual-time or profile-time clocks which measure CPU time in two different ways. These
5882 functions do not support access to the execution time of other processes.
- 5883 IBM's MVS operating system supports per-process and per-thread execution time clocks. It
5884 also supports limiting the execution time of a given process.
- 5885 Given all this existing practice, the working group considered that the POSIX.1 clocks and
5886 timers interface was appropriate to meet most of the requirements that realtime applications
5887 have for execution time clocks. Functions were added to get the CPU time clock IDs, and to
5888 allow/disallow the thread CPU-time clocks (in order to preserve the efficiency of some
5889 implementations of threads).
- 5890 • Clock Constants
- 5891 The definition of the manifest constants `CLOCK_PROCESS_CPUTIME_ID` and
5892 `CLOCK_THREAD_CPUTIME_ID` allows processes or threads, respectively, to access their
5893 own execution-time clocks. However, given a process or thread, access to its own execution-
5894 time clock is also possible if the clock ID of this clock is obtained through a call to
5895 *clock_getcpuclockid()* or *pthread_getcpuclockid()*. Therefore, these constants are not necessary
5896 and could be deleted to make the interface simpler. Their existence saves one system call in
5897 the first access to the CPU-time clock of each process or thread. The working group
5898 considered this issue and decided to leave the constants in IEEE Std 1003.1-2001 because they
5899 are closer to the POSIX.1b use of clock identifiers.
- 5900 • Library Implementations of Threads
- 5901 In library implementations of threads, kernel entities and library threads can coexist. In this
5902 case, if the CPU-time clocks are supported, most of the clock and timer functions will need to

5903 have two implementations: one in the thread library, and one in the system calls library. The
 5904 main difference between these two implementations is that the thread library
 5905 implementation will have to deal with clocks and timers that reside in the thread space,
 5906 while the kernel implementation will operate on timers and clocks that reside in kernel space.
 5907 In the library implementation, if the clock ID refers to a clock that resides in the kernel, a
 5908 kernel call will have to be made. The correct version of the function can be chosen by
 5909 specifying the appropriate order for the libraries during the link process.

5910 • History of Resolution Issues: Deletion of the *enable* Attribute

5911 In early proposals, consideration was given to inclusion of an attribute called *enable* for CPU-
 5912 time clocks. This would allow implementations to avoid the overhead of measuring
 5913 execution time for those processes or threads for which this measurement was not required.
 5914 However, this is unnecessary since processes are already required to measure execution time
 5915 by the POSIX.1 *times()* function. Consequently, the *enable* attribute is not present.

5916 Rationale Relating to Timeouts

5917 • Requirements for Timeouts

5918 Realtime systems which must operate reliably over extended periods without human
 5919 intervention are characteristic in embedded applications such as avionics, machine control,
 5920 and space exploration, as well as more mundane applications such as cable TV, security
 5921 systems, and plant automation. A multi-tasking paradigm, in which many independent
 5922 and/or cooperating software functions relinquish the processor(s) while waiting for a
 5923 specific stimulus, resource, condition, or operation completion, is very useful in producing
 5924 well engineered programs for such systems. For such systems to be robust and fault-tolerant,
 5925 expected occurrences that are unduly delayed or that never occur must be detected so that
 5926 appropriate recovery actions may be taken. This is difficult if there is no way for a task to
 5927 regain control of a processor once it has relinquished control (blocked) awaiting an
 5928 occurrence which, perhaps because of corrupted code, hardware malfunction, or latent
 5929 software bugs, will not happen when expected. Therefore, the common practice in realtime
 5930 operating systems is to provide a capability to time out such blocking services. Although
 5931 there are several methods to achieve this already defined by POSIX, none are as reliable or
 5932 efficient as initiating a timeout simultaneously with initiating a blocking service. This is
 5933 especially critical in hard-realtime embedded systems because the processors typically have
 5934 little time reserve, and allowed fault recovery times are measured in milliseconds rather than
 5935 seconds.

5936 The working group largely agreed that such timeouts were necessary and ought to become
 5937 part of IEEE Std 1003.1-2001, particularly vendors of realtime operating systems whose
 5938 customers had already expressed a strong need for timeouts. There was some resistance to
 5939 inclusion of timeouts in IEEE Std 1003.1-2001 because the desired effect, fault tolerance,
 5940 could, in theory, be achieved using existing facilities and alternative software designs, but
 5941 there was no compelling evidence that realtime system designers would embrace such
 5942 designs at the sacrifice of performance and/or simplicity.

5943 • Which Services should be Timed Out?

5944 Originally, the working group considered the prospect of providing timeouts on all blocking
 5945 services, including those currently existing in POSIX.1, POSIX.1b, and POSIX.1c, and future
 5946 interfaces to be defined by other working groups, as sort of a general policy. This was rather
 5947 quickly rejected because of the scope of such a change, and the fact that many of those
 5948 services would not normally be used in a realtime context. More traditional timesharing
 5949 solutions to timeout would suffice for most of the POSIX.1 interfaces, while others had
 5950 asynchronous alternatives which, while more complex to utilize, would be adequate for

5951 some realtime and all non-realtime applications.

5952 The list of potential candidates for timeouts was narrowed to the following for further
5953 consideration:

5954 — POSIX.1b

5955 — *sem_wait()*

5956 — *mq_receive()*

5957 — *mq_send()*

5958 — *lio_listio()*

5959 — *aio_suspend()*

5960 — *sigwait()* (timeout already implemented by *sigtimedwait()*)

5961 — POSIX.1c

5962 — *pthread_mutex_lock()*

5963 — *pthread_join()*

5964 — *pthread_cond_wait()* (timeout already implemented by *pthread_cond_timedwait()*)

5965 — POSIX.1

5966 — *read()*

5967 — *write()*

5968 After further review by the working group, the *lio_listio()*, *read()*, and *write()* functions (all
5969 forms of blocking synchronous I/O) were eliminated from the list because of the following:

5970 — Asynchronous alternatives exist

5971 — Timeouts can be implemented, albeit non-portably, in device drivers

5972 — A strong desire not to introduce modifications to POSIX.1 interfaces

5973 The working group ultimately rejected *pthread_join()* since both that interface and a timed
5974 variant of that interface are non-minimal and may be implemented as a function. See below
5975 for a library implementation of *pthread_join()*.

5976 Thus, there was a consensus among the working group members to add timeouts to 4 of the
5977 remaining 5 functions (the timeout for *aio_suspend()* was ultimately added directly to
5978 POSIX.1b, while the others were added by POSIX.1d). However, *pthread_mutex_lock()*
5979 remained contentious.

5980 Many feel that *pthread_mutex_lock()* falls into the same class as the other functions; that is, it
5981 is desirable to time out a mutex lock because a mutex may fail to be unlocked due to errant or
5982 corrupted code in a critical section (looping or branching outside of the unlock code), and
5983 therefore is equally in need of a reliable, simple, and efficient timeout. In fact, since mutexes
5984 are intended to guard small critical sections, most *pthread_mutex_lock()* calls would be
5985 expected to obtain the lock without blocking nor utilizing any kernel service, even in
5986 implementations of threads with global contention scope; the timeout alternative need only
5987 be considered after it is determined that the thread must block.

5988 Those opposed to timing out mutexes feel that the very simplicity of the mutex is
5989 compromised by adding a timeout semantic, and that to do so is senseless. They claim that if
5990 a timed mutex is really deemed useful by a particular application, then it can be constructed
5991 from the facilities already in POSIX.1b and POSIX.1c. The following two C-language library

5992 implementations of mutex locking with timeout represent the solutions offered (in both
 5993 implementations, the timeout parameter is specified as absolute time, not relative time as in
 5994 the proposed POSIX.1c interfaces).

5995 • Spinlock Implementation

```

5996     #include <pthread.h>
5997     #include <time.h>
5998     #include <errno.h>

5999     int pthread_mutex_timedlock(pthread_mutex_t *mutex,
6000                               const struct timespec *timeout)
6001     {
6002         struct timespec timenow;

6003         while (pthread_mutex_trylock(mutex) == EBUSY)
6004             {
6005                 clock_gettime(CLOCK_REALTIME, &timenow);
6006                 if (timespec_cmp(&timenow, timeout) >= 0)
6007                     {
6008                         return ETIMEDOUT;
6009                     }
6010                 pthread_yield();
6011             }
6012         return 0;
6013     }

```

6014 The Spinlock implementation is generally unsuitable for any application using priority-based
 6015 thread scheduling policies such as SCHED_FIFO or SCHED_RR, since the mutex could
 6016 currently be held by a thread of lower priority within the same allocation domain, but since
 6017 the waiting thread never blocks, only threads of equal or higher priority will ever run, and
 6018 the mutex cannot be unlocked. Setting priority inheritance or priority ceiling protocol on the
 6019 mutex does not solve this problem, since the priority of a mutex owning thread is only
 6020 boosted if higher priority threads are blocked waiting for the mutex; clearly not the case for
 6021 this spinlock.

6022 • Condition Wait Implementation

```

6023     #include <pthread.h>
6024     #include <time.h>
6025     #include <errno.h>

6026     struct timed_mutex
6027     {
6028         int locked;
6029         pthread_mutex_t mutex;
6030         pthread_cond_t cond;
6031     };
6032     typedef struct timed_mutex timed_mutex_t;

6033     int timed_mutex_lock(timed_mutex_t *tm,
6034                       const struct timespec *timeout)
6035     {
6036         int timedout=FALSE;
6037         int error_status;

```

```

6038     pthread_mutex_lock(&tm->mutex);
6039     while (tm->locked && !timedout)
6040     {
6041         if ((error_status=pthread_cond_timedwait(&tm->cond,
6042         &tm->mutex,
6043         timeout))!=0)
6044         {
6045             if (error_status==ETIMEDOUT) timedout = TRUE;
6046         }
6047     }
6048     if(timedout)
6049     {
6050         pthread_mutex_unlock(&tm->mutex);
6051         return ETIMEDOUT;
6052     }
6053     else
6054     {
6055         tm->locked = TRUE;
6056         pthread_mutex_unlock(&tm->mutex);
6057         return 0;
6058     }
6059 }
6060 void timed_mutex_unlock(timed_mutex_t *tm)
6061 {
6062     pthread_mutex_lock(&tm->mutex); / for case assignment not atomic /
6063     tm->locked = FALSE;
6064     pthread_mutex_unlock(&tm->mutex);
6065     pthread_cond_signal(&tm->cond);
6066 }

```

6067 The Condition Wait implementation effectively substitutes the *pthread_cond_timedwait()*
6068 function (which is currently timed out) for the desired *pthread_mutex_timedlock()*. Since waits
6069 on condition variables currently do not include protocols which avoid priority inversion, this
6070 method is generally unsuitable for realtime applications because it does not provide the same
6071 priority inversion protection as the untimed *pthread_mutex_lock()*. Also, for any given
6072 implementations of the current mutex and condition variable primitives, this library
6073 implementation has a performance cost at least 2.5 times that of the untimed
6074 *pthread_mutex_lock()* even in the case where the timed mutex is readily locked without
6075 blocking (the interfaces required for this case are shown in bold). Even in uniprocessors or
6076 where assignment is atomic, at least an additional *pthread_cond_signal()* is required.
6077 *pthread_mutex_timedlock()* could be implemented at effectively no performance penalty in
6078 this case because the timeout parameters need only be considered after it is determined that
6079 the mutex cannot be locked immediately.

6080 Thus it has not yet been shown that the full semantics of mutex locking with timeout can be
6081 efficiently and reliably achieved using existing interfaces. Even if the existence of an
6082 acceptable library implementation were proven, it is difficult to justify why the interface
6083 itself should not be made portable, especially considering approval for the other four
6084 timeouts.

```

6085     • Rationale for Library Implementation of pthread_timedjoin()
6086     Library implementation of pthread_timedjoin():
6087     /*
6088     * Construct a thread variety entirely from existing functions
6089     * with which a join can be done, allowing the join to time out.
6090     */
6091     #include <pthread.h>
6092     #include <time.h>
6093     struct timed_thread {
6094         pthread_t t;
6095         pthread_mutex_t m;
6096         int exiting;
6097         pthread_cond_t exit_c;
6098         void *(*start_routine)(void *arg);
6099         void *arg;
6100         void *status;
6101     };
6102     typedef struct timed_thread *timed_thread_t;
6103     static pthread_key_t timed_thread_key;
6104     static pthread_once_t timed_thread_once = PTHREAD_ONCE_INIT;
6105     static void timed_thread_init()
6106     {
6107         pthread_key_create(&timed_thread_key, NULL);
6108     }
6109     static void *timed_thread_start_routine(void *args)
6110     /*
6111     * Routine to establish thread-specific data value and run the actual
6112     * thread start routine which was supplied to timed_thread_create().
6113     */
6114     {
6115         timed_thread_t tt = (timed_thread_t) args;
6116         pthread_once(&timed_thread_once, timed_thread_init);
6117         pthread_setspecific(timed_thread_key, (void *)tt);
6118         timed_thread_exit((tt->start_routine)(tt->arg));
6119     }
6120     int timed_thread_create(timed_thread_t ttp, const pthread_attr_t *attr,
6121         void *(*start_routine)(void *), void *arg)
6122     /*
6123     * Allocate a thread which can be used with timed_thread_join().
6124     */
6125     {
6126         timed_thread_t tt;
6127         int result;
6128         tt = (timed_thread_t) malloc(sizeof(struct timed_thread));
6129         pthread_mutex_init(&tt->m, NULL);
6130         tt->exiting = FALSE;
6131         pthread_cond_init(&tt->exit_c, NULL);

```

```

6132         tt->start_routine = start_routine;
6133         tt->arg = arg;
6134         tt->status = NULL;

6135         if ((result = pthread_create(&tt->t, attr,
6136             timed_thread_start_routine, (void *)tt)) != 0) {
6137             free(tt);
6138             return result;
6139         }

6140         pthread_detach(tt->t);
6141         ttp = tt;
6142         return 0;
6143     }

6144     int timed_thread_join(timed_thread_t tt,
6145         struct timespec *timeout,
6146         void **status)
6147     {
6148         int result;

6149         pthread_mutex_lock(&tt->m);
6150         result = 0;
6151         /*
6152          * Wait until the thread announces that it is exiting,
6153          * or until timeout.
6154          */
6155         while (result == 0 && ! tt->exiting) {
6156             result = pthread_cond_timedwait(&tt->exit_c, &tt->m, timeout);
6157         }
6158         pthread_mutex_unlock(&tt->m);
6159         if (result == 0 && tt->exiting) {
6160             *status = tt->status;
6161             free((void *)tt);
6162             return result;
6163         }
6164         return result;
6165     }

6166     void timed_thread_exit(void *status)
6167     {
6168         timed_thread_t tt;
6169         void *specific;

6170         if ((specific=pthread_getspecific(timed_thread_key)) == NULL){
6171             /*
6172              * Handle cases which won't happen with correct usage.
6173              */
6174             pthread_exit( NULL);
6175         }
6176         tt = (timed_thread_t) specific;
6177         pthread_mutex_lock(&tt->m);
6178         /*
6179          * Tell a joiner that we're exiting.
6180          */

```

```

6181         tt->status = status;
6182         tt->exiting = TRUE;
6183         pthread_cond_signal(&tt->exit_c);
6184         pthread_mutex_unlock(&tt->m);
6185         /*
6186          * Call pthread_exit() to call destructors and really
6187          * exit the thread.
6188          */
6189         pthread_exit(NULL);
6190     }

```

6191 The *pthread_join()* C-language example shown above demonstrates that it is possible, using
6192 existing pthread facilities, to construct a variety of thread which allows for joining such a
6193 thread, but which allows the join operation to time out. It does this by using a
6194 *pthread_cond_timedwait()* to wait for the thread to exit. A **timed_thread_t** descriptor structure
6195 is used to pass parameters from the creating thread to the created thread, and from the
6196 exiting thread to the joining thread. This implementation is roughly equivalent to what a
6197 normal *pthread_join()* implementation would do, with the single change being that
6198 *pthread_cond_timedwait()* is used in place of a simple *pthread_cond_wait()*.

6199 Since it is possible to implement such a facility entirely from existing pthread interfaces, and
6200 with roughly equal efficiency and complexity to an implementation which would be
6201 provided directly by a pthreads implementation, it was the consensus of the working group
6202 members that any *pthread_timedjoin()* facility would be unnecessary, and should not be
6203 provided.

6204 • Form of the Timeout Interfaces

6205 The working group considered a number of alternative ways to add timeouts to blocking
6206 services. At first, a system interface which would specify a one-shot or persistent timeout to
6207 be applied to subsequent blocking services invoked by the calling process or thread was
6208 considered because it allowed all blocking services to be timed out in a uniform manner with
6209 a single additional interface; this was rather quickly rejected because it could easily result in
6210 the wrong services being timed out.

6211 It was suggested that a timeout value might be specified as an attribute of the object
6212 (semaphore, mutex, message queue, and so on), but there was no consensus on this, either on
6213 a case-by-case basis or for all timeouts.

6214 Looking at the two existing timeouts for blocking services indicates that the working group
6215 members favor a separate interface for the timed version of a function. However,
6216 *pthread_cond_timedwait()* utilizes an absolute timeout value while *sigtimedwait()* uses a
6217 relative timeout value. The working group members agreed that relative timeout values are
6218 appropriate where the timeout mechanism's primary use was to deal with an unexpected or
6219 error situation, but they are inappropriate when the timeout must expire at a particular time,
6220 or before a specific deadline. For the timeouts being introduced in IEEE Std 1003.1-2001, the
6221 working group considered allowing both relative and absolute timeouts as is done with
6222 POSIX.1b timers, but ultimately favored the simpler absolute timeout form.

6223 An absolute time measure can be easily implemented on top of an interface that specifies
6224 relative time, by reading the clock, calculating the difference between the current time and
6225 the desired wake-up time, and issuing a relative timeout call. But there is a race condition
6226 with this approach because the thread could be preempted after reading the clock, but before
6227 making the timed-out call; in this case, the thread would be awakened later than it should
6228 and, thus, if the wake-up time represented a deadline, it would miss it.

6229 There is also a race condition when trying to build a relative timeout on top of an interface
6230 that specifies absolute timeouts. In this case, the clock would have to be read to calculate the
6231 absolute wake-up time as the sum of the current time plus the relative timeout interval. In
6232 this case, if the thread is preempted after reading the clock but before making the timed-out
6233 call, the thread would be awakened earlier than desired.

6234 But the race condition with the absolute timeouts interface is not as bad as the one that
6235 happens with the relative timeout interface, because there are simple workarounds. For the
6236 absolute timeouts interface, if the timing requirement is a deadline, the deadline can still be
6237 met because the thread woke up earlier than the deadline. If the timeout is just used as an
6238 error recovery mechanism, the precision of timing is not really important. If the timing
6239 requirement is that between actions A and B a minimum interval of time must elapse, the
6240 absolute timeout interface can be safely used by reading the clock after action A has been
6241 started. It could be argued that, since the call with the absolute timeout is atomic from the
6242 application point of view, it is not possible to read the clock after action A, if this action is
6243 part of the timed-out call. But looking at the nature of the calls for which timeouts are
6244 specified (locking a mutex, waiting for a semaphore, waiting for a message, or waiting until
6245 there is space in a message queue), the timeouts that an application would build on these
6246 actions would not be triggered by these actions themselves, but by some other external
6247 action. For example, if waiting for a message to arrive to a message queue, and waiting for at
6248 least 20 milliseconds, this time interval would start to be counted from some event that
6249 would trigger both the action that produces the message, as well as the action that waits for
6250 the message to arrive, and not by the wait-for-message operation itself. In this case, the
6251 workaround proposed above could be used.

6252 For these reasons, the absolute timeout is preferred over the relative timeout interface.

6253 **B.2.9 Threads**

6254 Threads will normally be more expensive than subroutines (or functions, routines, and so on) if
6255 specialized hardware support is not provided. Nevertheless, threads should be sufficiently
6256 efficient to encourage their use as a medium to fine-grained structuring mechanism for
6257 parallelism in an application. Structuring an application using threads then allows it to take
6258 immediate advantage of any underlying parallelism available in the host environment. This
6259 means implementors are encouraged to optimize for fast execution at the possible expense of
6260 efficient utilization of storage. For example, a common thread creation technique is to cache
6261 appropriate thread data structures. That is, rather than releasing system resources, the
6262 implementation retains these resources and reuses them when the program next asks to create a
6263 new thread. If this reuse of thread resources is to be possible, there has to be very little unique
6264 state associated with each thread, because any such state has to be reset when the thread is
6265 reused.

6266 **Thread Creation Attributes**

6267 Attributes objects are provided for threads, mutexes, and condition variables as a mechanism to
6268 support probable future standardization in these areas without requiring that the interface itself
6269 be changed.

6270 Attributes objects provide clean isolation of the configurable aspects of threads. For example,
6271 “stack size” is an important attribute of a thread, but it cannot be expressed portably. When
6272 porting a threaded program, stack sizes often need to be adjusted. The use of attributes objects
6273 can help by allowing the changes to be isolated in a single place, rather than being spread across
6274 every instance of thread creation.

6275 Attributes objects can be used to set up *classes* of threads with similar attributes; for example,
6276 “threads with large stacks and high priority” or “threads with minimal stacks”. These classes
6277 can be defined in a single place and then referenced wherever threads need to be created.
6278 Changes to “class” decisions become straightforward, and detailed analysis of each
6279 *pthread_create()* call is not required.

6280 The attributes objects are defined as opaque types as an aid to extensibility. If these objects had
6281 been specified as structures, adding new attributes would force recompilation of all multi-
6282 threaded programs when the attributes objects are extended; this might not be possible if
6283 different program components were supplied by different vendors.

6284 Additionally, opaque attributes objects present opportunities for improving performance.
6285 Argument validity can be checked once when attributes are set, rather than each time a thread is
6286 created. Implementations will often need to cache kernel objects that are expensive to create.
6287 Opaque attributes objects provide an efficient mechanism to detect when cached objects become
6288 invalid due to attribute changes.

6289 Because assignment is not necessarily defined on a given opaque type, implementation-defined
6290 default values cannot be defined in a portable way. The solution to this problem is to allow
6291 attribute objects to be initialized dynamically by attributes object initialization functions, so that
6292 default values can be supplied automatically by the implementation.

6293 The following proposal was provided as a suggested alternative to the supplied attributes:

- 6294 1. Maintain the style of passing a parameter formed by the bitwise-inclusive OR of flags to
6295 the initialization routines (*pthread_create()*, *pthread_mutex_init()*, *pthread_cond_init()*). The
6296 parameter containing the flags should be an opaque type for extensibility. If no flags are
6297 set in the parameter, then the objects are created with default characteristics. An
6298 implementation may specify implementation-defined flag values and associated behavior.
- 6299 2. If further specialization of mutexes and condition variables is necessary, implementations
6300 may specify additional procedures that operate on the **pthread_mutex_t** and
6301 **pthread_cond_t** objects (instead of on attributes objects).

6302 The difficulties with this solution are:

- 6303 1. A bitmask is not opaque if bits have to be set into bit-vector attributes objects using
6304 explicitly-coded bitwise-inclusive OR operations. If the set of options exceeds an **int**,
6305 application programmers need to know the location of each bit. If bits are set or read by
6306 encapsulation (that is, *get*()* or *set*()* functions), then the bitmask is merely an
6307 implementation of attributes objects as currently defined and should not be exposed to the
6308 programmer.
- 6309 2. Many attributes are not Boolean or very small integral values. For example, scheduling
6310 policy may be placed in 3 bits or 4 bits, but priority requires 5 bits or more, thereby taking
6311 up at least 8 bits out of a possible 16 bits on machines with 16-bit integers. Because of this,
6312 the bitmask can only reasonably control whether particular attributes are set or not, and it
6313 cannot serve as the repository of the value itself. The value needs to be specified as a
6314 function parameter (which is non-extensible), or by setting a structure field (which is non-
6315 opaque), or by *get*()* and *set*()* functions (making the bitmask a redundant addition to the
6316 attributes objects).

6317 Stack size is defined as an optional attribute because the very notion of a stack is inherently
6318 machine-dependent. Some implementations may not be able to change the size of the stack, for
6319 example, and others may not need to because stack pages may be discontinuous and can be
6320 allocated and released on demand.

6321 The attribute mechanism has been designed in large measure for extensibility. Future extensions
6322 to the attribute mechanism or to any attributes object defined in IEEE Std 1003.1-2001 have to be
6323 done with care so as not to affect binary-compatibility.

6324 Attribute objects, even if allocated by means of dynamic allocation functions such as *malloc()*,
6325 may have their size fixed at compile time. This means, for example, a *pthread_create()* in an
6326 implementation with extensions to the **pthread_attr_t** cannot look beyond the area that the
6327 binary application assumes is valid. This suggests that implementations should maintain a size
6328 field in the attributes object, as well as possibly version information, if extensions in different
6329 directions (possibly by different vendors) are to be accommodated.

6330 Thread Implementation Models

6331 There are various thread implementation models. At one end of the spectrum is the “library-
6332 thread model”. In such a model, the threads of a process are not visible to the operating system
6333 kernel, and the threads are not kernel-scheduled entities. The process is the only kernel-
6334 scheduled entity. The process is scheduled onto the processor by the kernel according to the
6335 scheduling attributes of the process. The threads are scheduled onto the single kernel-scheduled
6336 entity (the process) by the runtime library according to the scheduling attributes of the threads.
6337 A problem with this model is that it constrains concurrency. Since there is only one kernel-
6338 scheduled entity (namely, the process), only one thread per process can execute at a time. If the
6339 thread that is executing blocks on I/O, then the whole process blocks.

6340 At the other end of the spectrum is the “kernel-thread model”. In this model, all threads are
6341 visible to the operating system kernel. Thus, all threads are kernel-scheduled entities, and all
6342 threads can concurrently execute. The threads are scheduled onto processors by the kernel
6343 according to the scheduling attributes of the threads. The drawback to this model is that the
6344 creation and management of the threads entails operating system calls, as opposed to subroutine
6345 calls, which makes kernel threads heavier weight than library threads.

6346 Hybrids of these two models are common. A hybrid model offers the speed of library threads
6347 and the concurrency of kernel threads. In hybrid models, a process has some (relatively small)
6348 number of kernel scheduled entities associated with it. It also has a potentially much larger
6349 number of library threads associated with it. Some library threads may be bound to kernel-
6350 scheduled entities, while the other library threads are multiplexed onto the remaining kernel-
6351 scheduled entities. There are two levels of thread scheduling:

- 6352 1. The runtime library manages the scheduling of (unbound) library threads onto kernel-
6353 scheduled entities.
- 6354 2. The kernel manages the scheduling of kernel-scheduled entities onto processors.

6355 For this reason, a hybrid model is referred to as a two-level threads scheduling model. In this
6356 model, the process can have multiple concurrently executing threads; specifically, it can have as
6357 many concurrently executing threads as it has kernel-scheduled entities.

6358 Thread-Specific Data

6359 Many applications require that a certain amount of context be maintained on a per-thread basis
6360 across procedure calls. A common example is a multi-threaded library routine that allocates
6361 resources from a common pool and maintains an active resource list for each thread. The
6362 thread-specific data interface provided to meet these needs may be viewed as a two-dimensional
6363 array of values with keys serving as the row index and thread IDs as the column index (although
6364 the implementation need not work this way).

- 6365 • Models

6366 Three possible thread-specific data models were considered:

6367 1. No Explicit Support

6368 A standard thread-specific data interface is not strictly necessary to support
 6369 applications that require per-thread context. One could, for example, provide a hash
 6370 function that converted a **pthread_t** into an integer value that could then be used to
 6371 index into a global array of per-thread data pointers. This hash function, in conjunction
 6372 with *pthread_self()*, would be all the interface required to support a mechanism of this
 6373 sort. Unfortunately, this technique is cumbersome. It can lead to duplicated code as
 6374 each set of cooperating modules implements their own per-thread data management
 6375 schemes.

6376 2. Single (**void ***) Pointer

6377 Another technique would be to provide a single word of per-thread storage and a pair
 6378 of functions to fetch and store the value of this word. The word could then hold a
 6379 pointer to a block of per-thread memory. The allocation, partitioning, and general use
 6380 of this memory would be entirely up to the application. Although this method is not as
 6381 problematic as technique 1, it suffers from interoperability problems. For example, all
 6382 modules using the per-thread pointer would have to agree on a common usage
 6383 protocol.

6384 3. Key/Value Mechanism

6385 This method associates an opaque key (for example, stored in a variable of type
 6386 **pthread_key_t**) with each per-thread datum. These keys play the role of identifiers for
 6387 per-thread data. This technique is the most generic and avoids the problems noted
 6388 above, albeit at the cost of some complexity.

6389 The primary advantage of the third model is its information hiding properties. Modules
 6390 using this model are free to create and use their own key(s) independent of all other such
 6391 usage, whereas the other models require that all modules that use thread-specific context
 6392 explicitly cooperate with all other such modules. The data-independence provided by the
 6393 third model is worth the additional interface.

6394 • Requirements

6395 It is important that it be possible to implement the thread-specific data interface without the
 6396 use of thread private memory. To do otherwise would increase the weight of each thread,
 6397 thereby limiting the range of applications for which the threads interfaces provided by
 6398 IEEE Std 1003.1-2001 is appropriate.

6399 The values that one binds to the key via *pthread_setspecific()* may, in fact, be pointers to
 6400 shared storage locations available to all threads. It is only the key/value bindings that are
 6401 maintained on a per-thread basis, and these can be kept in any portion of the address space
 6402 that is reserved for use by the calling thread (for example, on the stack). Thus, no per-thread
 6403 MMU state is required to implement the interface. On the other hand, there is nothing in the
 6404 interface specification to preclude the use of a per-thread MMU state if it is available (for
 6405 example, the key values returned by *pthread_key_create()* could be thread private memory
 6406 addresses).

6407 • Standardization Issues

6408 Thread-specific data is a requirement for a usable thread interface. The binding described in
 6409 this section provides a portable thread-specific data mechanism for languages that do not
 6410 directly support a thread-specific storage class. A binding to IEEE Std 1003.1-2001 for a
 6411 language that does include such a storage class need not provide this specific interface.

6412 If a language were to include the notion of thread-specific storage, it would be desirable (but
 6413 *not* required) to provide an implementation of the pthreads thread-specific data interface
 6414 based on the language feature. For example, assume that a compiler for a C-like language
 6415 supports a *private* storage class that provides thread-specific storage. Something similar to
 6416 the following macros might be used to effect a compatible implementation:

```
6417     #define pthread_key_t           private void *
6418     #define pthread_key_create(key) /* no-op */
6419     #define pthread_setspecific(key,value) (key)=(value)
6420     #define pthread_getspecific(key)     (key)
```

6421 **Note:** For the sake of clarity, this example ignores destructor functions. A correct implementation
 6422 would have to support them.

6423 Barriers

6424 • Background

6425 Barriers are typically used in parallel DO/FOR loops to ensure that all threads have reached
 6426 a particular stage in a parallel computation before allowing any to proceed to the next stage.
 6427 Highly efficient implementation is possible on machines which support a “Fetch and Add”
 6428 operation as described in the referenced Almasi and Gottlieb (1989).

6429 The use of return value PTHREAD_BARRIER_SERIAL_THREAD is shown in the following
 6430 example:

```
6431     if ( (status=pthread_barrier_wait(&barrier)) ==
6432         PTHREAD_BARRIER_SERIAL_THREAD) {
6433         ...serial section
6434     }
6435         else if (status != 0) {
6436             ...error processing
6437         }
6438     status=pthread_barrier_wait(&barrier);
6439     ...
```

6440 This behavior allows a serial section of code to be executed by one thread as soon as all
 6441 threads reach the first barrier. The second barrier prevents the other threads from proceeding
 6442 until the serial section being executed by the one thread has completed.

6443 Although barriers can be implemented with mutexes and condition variables, the referenced
 6444 Almasi and Gottlieb (1989) provides ample illustration that such implementations are
 6445 significantly less efficient than is possible. While the relative efficiency of barriers may well
 6446 vary by implementation, it is important that they be recognized in the IEEE Std 1003.1-2001
 6447 to facilitate applications portability while providing the necessary freedom to implementors.

6448 • Lack of Timeout Feature

6449 Alternate versions of most blocking routines have been provided to support watchdog
 6450 timeouts. No alternate interface of this sort has been provided for barrier waits for the
 6451 following reasons:

- 6452 • Multiple threads may use different timeout values, some of which may be indefinite. It is
 6453 not clear which threads should break through the barrier with a timeout error if and when
 6454 these timeouts expire.
- 6455 • The barrier may become unusable once a thread breaks out of a *pthread_barrier_wait()*
 6456 with a timeout error. There is, in general, no way to guarantee the consistency of a

6457 barrier's internal data structures once a thread has timed out of a `pthread_barrier_wait()`.
6458 Even the inclusion of a special barrier reinitialization function would not help much since
6459 it is not clear how this function would affect the behavior of threads that reach the barrier
6460 between the original timeout and the call to the reinitialization function.

6461 Spin Locks

6462 • Background

6463 Spin locks represent an extremely low-level synchronization mechanism suitable primarily
6464 for use on shared memory multi-processors. It is typically an atomically modified Boolean
6465 value that is set to one when the lock is held and to zero when the lock is freed.

6466 When a caller requests a spin lock that is already held, it typically spins in a loop testing
6467 whether the lock has become available. Such spinning wastes processor cycles so the lock
6468 should only be held for short durations and not across sleep/block operations. Callers should
6469 unlock spin locks before calling sleep operations.

6470 Spin locks are available on a variety of systems. The functions included in
6471 IEEE Std 1003.1-2001 are an attempt to standardize that existing practice.

6472 • Lack of Timeout Feature

6473 Alternate versions of most blocking routines have been provided to support watchdog
6474 timeouts. No alternate interface of this sort has been provided for spin locks for the following
6475 reasons:

- 6476 • It is impossible to determine appropriate timeout intervals for spin locks in a portable
6477 manner. The amount of time one can expect to spend spin-waiting is inversely
6478 proportional to the degree of parallelism provided by the system.

6479 It can vary from a few cycles when each competing thread is running on its own
6480 processor, to an indefinite amount of time when all threads are multiplexed on a single
6481 processor (which is why spin locking is not advisable on uniprocessors).

- 6482 • When used properly, the amount of time the calling thread spends waiting on a spin lock
6483 should be considerably less than the time required to set up a corresponding watchdog
6484 timer. Since the primary purpose of spin locks is to provide a low-overhead
6485 synchronization mechanism for multi-processors, the overhead of a timeout mechanism
6486 was deemed unacceptable.

6487 It was also suggested that an additional *count* argument be provided (on the
6488 `pthread_spin_lock()` call) in lieu of a true timeout so that a spin lock call could fail gracefully if
6489 it was unable to apply the lock after *count* attempts. This idea was rejected because it is not
6490 existing practice. Furthermore, the same effect can be obtained with `pthread_spin_trylock()`,
6491 as illustrated below:

```

6492         int n = MAX_SPIN;
6493         while ( --n >= 0 )
6494         {
6495             if ( !pthread_spin_try_lock(...) )
6496                 break;
6497         }
6498         if ( n >= 0 )
6499         {
6500             /* Successfully acquired the lock */
6501         }
6502         else
6503         {
6504             /* Unable to acquire the lock */
6505         }

```

6506 • *process-shared* Attribute

6507 The initialization functions associated with most POSIX synchronization objects (for
6508 example, mutexes, barriers, and read-write locks) take an attributes object with a *process-*
6509 *shared* attribute that specifies whether or not the object is to be shared across processes. In the
6510 draft corresponding to the first balloting round, two separate initialization functions are
6511 provided for spin locks, however: one for spin locks that were to be shared across processes
6512 (*spin_init()*), and one for locks that were only used by multiple threads within a single
6513 process (*pthread_spin_init()*). This was done so as to keep the overhead associated with spin
6514 waiting to an absolute minimum. However, the balloting group requested that, since the
6515 overhead associated to a bit check was small, spin locks should be consistent with the rest of
6516 the synchronization primitives, and thus the *process-shared* attribute was introduced for spin
6517 locks.

6518 • Spin Locks *versus* Mutexes

6519 It has been suggested that mutexes are an adequate synchronization mechanism and spin
6520 locks are not necessary. Locking mechanisms typically must trade off the processor resources
6521 consumed while setting up to block the thread and the processor resources consumed by the
6522 thread while it is blocked. Spin locks require very little resources to set up the blocking of a
6523 thread. Existing practice is to simply loop, repeating the atomic locking operation until the
6524 lock is available. While the resources consumed to set up blocking of the thread are low, the
6525 thread continues to consume processor resources while it is waiting.

6526 On the other hand, mutexes may be implemented such that the processor resources
6527 consumed to block the thread are large relative to a spin lock. After detecting that the mutex
6528 lock is not available, the thread must alter its scheduling state, add itself to a set of waiting
6529 threads, and, when the lock becomes available again, undo all of this before taking over
6530 ownership of the mutex. However, while a thread is blocked by a mutex, no processor
6531 resources are consumed.

6532 Therefore, spin locks and mutexes may be implemented to have different characteristics.
6533 Spin locks may have lower overall overhead for very short-term blocking, and mutexes may
6534 have lower overall overhead when a thread will be blocked for longer periods of time. The
6535 presence of both interfaces allows implementations with these two different characteristics,
6536 both of which may be useful to a particular application.

6537 It has also been suggested that applications can build their own spin locks from the
6538 *pthread_mutex_trylock()* function:

```
6539     while (pthread_mutex_trylock(&mutex));
```

6540 The apparent simplicity of this construct is somewhat deceiving, however. While the actual
6541 wait is quite efficient, various guarantees on the integrity of mutex objects (for example,
6542 priority inheritance rules) may add overhead to the successful path of the trylock operation
6543 that is not required of spin locks. One could, of course, add an attribute to the mutex to
6544 bypass such overhead, but the very act of finding and testing this attribute represents more
6545 overhead than is found in the typical spin lock.

6546 The need to hold spin lock overhead to an absolute minimum also makes it impossible to
6547 provide guarantees against starvation similar to those provided for mutexes or read-write
6548 locks. The overhead required to implement such guarantees (for example, disabling
6549 preemption before spinning) may well exceed the overhead of the spin wait itself by many
6550 orders of magnitude. If a “safe” spin wait seems desirable, it can always be provided (albeit
6551 at some performance cost) via appropriate mutex attributes.

6552 XSI Supported Functions

6553 On XSI-conformant systems, the following symbolic constants are always defined:

```
6554     _POSIX_READER_WRITER_LOCKS
6555     _POSIX_THREAD_ATTR_STACKADDR
6556     _POSIX_THREAD_ATTR_STACKSIZE
6557     _POSIX_THREAD_PROCESS_SHARED
6558     _POSIX_THREADS
```

6559 Therefore, the following threads functions are always supported:

6560	<i>pthread_atfork()</i>	<i>pthread_key_create()</i>
6561	<i>pthread_attr_destroy()</i>	<i>pthread_key_delete()</i>
6562	<i>pthread_attr_getdetachstate()</i>	<i>pthread_kill()</i>
6563	<i>pthread_attr_getguardsize()</i>	<i>pthread_mutex_destroy()</i>
6564	<i>pthread_attr_getschedparam()</i>	<i>pthread_mutex_init()</i>
6565	<i>pthread_attr_getstack()</i>	<i>pthread_mutex_lock()</i>
6566	<i>pthread_attr_getstackaddr()</i>	<i>pthread_mutex_trylock()</i>
6567	<i>pthread_attr_getstacksize()</i>	<i>pthread_mutex_unlock()</i>
6568	<i>pthread_attr_init()</i>	<i>pthread_mutexattr_destroy()</i>
6569	<i>pthread_attr_setdetachstate()</i>	<i>pthread_mutexattr_getpshared()</i>
6570	<i>pthread_attr_setguardsize()</i>	<i>pthread_mutexattr_gettype()</i>
6571	<i>pthread_attr_setschedparam()</i>	<i>pthread_mutexattr_init()</i>
6572	<i>pthread_attr_setstack()</i>	<i>pthread_mutexattr_setpshared()</i>
6573	<i>pthread_attr_setstackaddr()</i>	<i>pthread_mutexattr_settype()</i>
6574	<i>pthread_attr_setstacksize()</i>	<i>pthread_once()</i>
6575	<i>pthread_cancel()</i>	<i>pthread_rwlock_destroy()</i>
6576	<i>pthread_cleanup_pop()</i>	<i>pthread_rwlock_init()</i>
6577	<i>pthread_cleanup_push()</i>	<i>pthread_rwlock_rdlock()</i>
6578	<i>pthread_cond_broadcast()</i>	<i>pthread_rwlock_tryrdlock()</i>
6579	<i>pthread_cond_destroy()</i>	<i>pthread_rwlock_trywrlock()</i>
6580	<i>pthread_cond_init()</i>	<i>pthread_rwlock_unlock()</i>
6581	<i>pthread_cond_signal()</i>	<i>pthread_rwlock_wrlock()</i>
6582	<i>pthread_cond_timedwait()</i>	<i>pthread_rwlockattr_destroy()</i>
6583	<i>pthread_cond_wait()</i>	<i>pthread_rwlockattr_getpshared()</i>
6584	<i>pthread_condattr_destroy()</i>	<i>pthread_rwlockattr_init()</i>

6585	<i>pthread_condattr_getpshared()</i>	<i>pthread_rwlockattr_setpshared()</i>
6586	<i>pthread_condattr_init()</i>	<i>pthread_self()</i>
6587	<i>pthread_condattr_setpshared()</i>	<i>pthread_setcancelstate()</i>
6588	<i>pthread_create()</i>	<i>pthread_setcanceltype()</i>
6589	<i>pthread_detach()</i>	<i>pthread_setconcurrency()</i>
6590	<i>pthread_equal()</i>	<i>pthread_setspecific()</i>
6591	<i>pthread_exit()</i>	<i>pthread_sigmask()</i>
6592	<i>pthread_getconcurrency()</i>	<i>pthread_testcancel()</i>
6593	<i>pthread_getspecific()</i>	<i>sigwait()</i>
6594	<i>pthread_join()</i>	

6595 On XSI-conformant systems, the symbolic constant `_POSIX_THREAD_SAFE_FUNCTIONS` is
6596 always defined. Therefore, the following functions are always supported:

6597	<i>asctime_r()</i>	<i>getpwuid_r()</i>
6598	<i>ctime_r()</i>	<i>gmtime_r()</i>
6599	<i>flockfile()</i>	<i>localtime_r()</i>
6600	<i>ftrylockfile()</i>	<i>putc_unlocked()</i>
6601	<i>funlockfile()</i>	<i>putchar_unlocked()</i>
6602	<i>getc_unlocked()</i>	<i>rand_r()</i>
6603	<i>getchar_unlocked()</i>	<i>readdir_r()</i>
6604	<i>getgrgid_r()</i>	<i>strerror_r()</i>
6605	<i>getgrnam_r()</i>	<i>strtok_r()</i>
6606	<i>getpwnam_r()</i>	

6607 The following threads functions are only supported on XSI-conformant systems if the Realtime
6608 Threads Option Group is supported :

6609	<i>pthread_attr_getinheritsched()</i>	<i>pthread_mutex_getprioceiling()</i>
6610	<i>pthread_attr_getschedpolicy()</i>	<i>pthread_mutex_setprioceiling()</i>
6611	<i>pthread_attr_getscope()</i>	<i>pthread_mutexattr_getprioceiling()</i>
6612	<i>pthread_attr_setinheritsched()</i>	<i>pthread_mutexattr_getprotocol()</i>
6613	<i>pthread_attr_setschedpolicy()</i>	<i>pthread_mutexattr_setprioceiling()</i>
6614	<i>pthread_attr_setscope()</i>	<i>pthread_mutexattr_setprotocol()</i>
6615	<i>pthread_getschedparam()</i>	<i>pthread_setschedparam()</i>

6616 XSI Threads Extensions

6617 The following XSI extensions to POSIX.1c are now supported in IEEE Std 1003.1-2001 as part of
6618 the alignment with the Single UNIX Specification:

- 6619 • Extended mutex attribute types
- 6620 • Read-write locks and attributes (also introduced by the IEEE Std 1003.1j-2000 amendment)
- 6621 • Thread concurrency level
- 6622 • Thread stack guard size
- 6623 • Parallel I/O

6624 A total of 19 new functions were added.

6625 These extensions carefully follow the threads programming model specified in POSIX.1c. As
6626 with POSIX.1c, all the new functions return zero if successful; otherwise, an error number is

6627 returned to indicate the error.

6628 The concept of attribute objects was introduced in POSIX.1c to allow implementations to extend
6629 IEEE Std 1003.1-2001 without changing the existing interfaces. Attribute objects were defined for
6630 threads, mutexes, and condition variables. Attributes objects are defined as implementation-
6631 defined opaque types to aid extensibility, and functions are defined to allow attributes to be set
6632 or retrieved. This model has been followed when adding the new type attribute of
6633 **pthread_mutexattr_t** or the new read-write lock attributes object **pthread_rwlockattr_t**.

6634 • Extended Mutex Attributes

6635 POSIX.1c defines a mutex attributes object as an implementation-defined opaque object of
6636 type **pthread_mutexattr_t**, and specifies a number of attributes which this object must have
6637 and a number of functions which manipulate these attributes. These attributes include
6638 *detachstate*, *inheritsched*, *schedparm*, *schedpolicy*, *contentionscope*, *stackaddr*, and *stacksize*.

6639 The System Interfaces volume of IEEE Std 1003.1-2001 specifies another mutex attribute
6640 called *type*. The *type* attribute allows applications to specify the behavior of mutex locking
6641 operations in situations where POSIX.1c behavior is undefined. The OSF DCE threads
6642 implementation, based on Draft 4 of POSIX.1c, specified a similar attribute. Note that the
6643 names of the attributes have changed somewhat from the OSF DCE threads implementation.

6644 The System Interfaces volume of IEEE Std 1003.1-2001 also extends the specification of the
6645 following POSIX.1c functions which manipulate mutexes:

6646 *pthread_mutex_lock()*
6647 *pthread_mutex_trylock()*
6648 *pthread_mutex_unlock()*

6649 to take account of the new mutex attribute type and to specify behavior which was declared
6650 as undefined in POSIX.1c. How a calling thread acquires or releases a mutex now depends
6651 upon the mutex *type* attribute.

6652 The *type* attribute can have the following values:

6653 PTHREAD_MUTEX_NORMAL
6654 Basic mutex with no specific error checking built in. Does not report a deadlock error.

6655 PTHREAD_MUTEX_RECURSIVE
6656 Allows any thread to recursively lock a mutex. The mutex must be unlocked an equal
6657 number of times to release the mutex.

6658 PTHREAD_MUTEX_ERRORCHECK
6659 Detects and reports simple usage errors; that is, an attempt to unlock a mutex that is not
6660 locked by the calling thread or that is not locked at all, or an attempt to relock a mutex
6661 the thread already owns.

6662 PTHREAD_MUTEX_DEFAULT
6663 The default mutex type. May be mapped to any of the above mutex types or may be an
6664 implementation-defined type.

6665 *Normal* mutexes do not detect deadlock conditions; for example, a thread will hang if it tries
6666 to relock a normal mutex that it already owns. Attempting to unlock a mutex locked by
6667 another thread, or unlocking an unlocked mutex, results in undefined behavior. Normal
6668 mutexes will usually be the fastest type of mutex available on a platform but provide the
6669 least error checking.

6670 *Recursive* mutexes are useful for converting old code where it is difficult to establish clear
6671 boundaries of synchronization. A thread can relock a recursive mutex without first unlocking

6672 it. The relocking deadlock which can occur with normal mutexes cannot occur with this type
6673 of mutex. However, multiple locks of a recursive mutex require the same number of unlocks
6674 to release the mutex before another thread can acquire the mutex. Furthermore, this type of
6675 mutex maintains the concept of an owner. Thus, a thread attempting to unlock a recursive
6676 mutex which another thread has locked returns with an error. A thread attempting to unlock
6677 a recursive mutex that is not locked returns with an error. Never use a recursive mutex with
6678 condition variables because the implicit unlock performed by *pthread_cond_wait()* or
6679 *pthread_cond_timedwait()* will not actually release the mutex if it had been locked multiple
6680 times.

6681 *Errorcheck* mutexes provide error checking and are useful primarily as a debugging aid. A
6682 thread attempting to relock an errorcheck mutex without first unlocking it returns with an
6683 error. Again, this type of mutex maintains the concept of an owner. Thus, a thread
6684 attempting to unlock an errorcheck mutex which another thread has locked returns with an
6685 error. A thread attempting to unlock an errorcheck mutex that is not locked also returns with
6686 an error. It should be noted that errorcheck mutexes will almost always be much slower than
6687 normal mutexes due to the extra state checks performed.

6688 The default mutex type provides implementation-defined error checking. The default mutex
6689 may be mapped to one of the other defined types or may be something entirely different.
6690 This enables each vendor to provide the mutex semantics which the vendor feels will be
6691 most useful to their target users. Most vendors will probably choose to make normal
6692 mutexes the default so as to give applications the benefit of the fastest type of mutexes
6693 available on their platform. Check your implementation's documentation.

6694 An application developer can use any of the mutex types almost interchangeably as long as
6695 the application does not depend upon the implementation detecting (or failing to detect) any
6696 particular errors. Note that a recursive mutex can be used with condition variable waits as
6697 long as the application never recursively locks the mutex.

6698 Two functions are provided for manipulating the *type* attribute of a mutex attributes object.
6699 This attribute is set or returned in the *type* parameter of these functions. The
6700 *pthread_mutexattr_settype()* function is used to set a specific type value while
6701 *pthread_mutexattr_gettype()* is used to return the type of the mutex. Setting the *type* attribute
6702 of a mutex attributes object affects only mutexes initialized using that mutex attributes
6703 object. Changing the *type* attribute does not affect mutexes previously initialized using that
6704 mutex attributes object.

6705 • Read-Write Locks and Attributes

6706 The read-write locks introduced have been harmonized with those in IEEE Std 1003.1j-2000;
6707 see also Section B.2.9.6 (on page 176).

6708 Read-write locks (also known as reader-writer locks) allow a thread to exclusively lock some
6709 shared data while updating that data, or allow any number of threads to have simultaneous
6710 read-only access to the data.

6711 Unlike a mutex, a read-write lock distinguishes between reading data and writing data. A
6712 mutex excludes all other threads. A read-write lock allows other threads access to the data,
6713 providing no thread is modifying the data. Thus, a read-write lock is less primitive than
6714 either a mutex-condition variable pair or a semaphore.

6715 Application developers should consider using a read-write lock rather than a mutex to
6716 protect data that is frequently referenced but seldom modified. Most threads (readers) will be
6717 able to read the data without waiting and will only have to block when some other thread (a
6718 writer) is in the process of modifying the data. Conversely a thread that wants to change the
6719 data is forced to wait until there are no readers. This type of lock is often used to facilitate

6720 parallel access to data on multi-processor platforms or to avoid context switches on single
6721 processor platforms where multiple threads access the same data.

6722 If a read-write lock becomes unlocked and there are multiple threads waiting to acquire the
6723 write lock, the implementation's scheduling policy determines which thread acquires the
6724 read-write lock for writing. If there are multiple threads blocked on a read-write lock for both
6725 read locks and write locks, it is unspecified whether the readers or a writer acquire the lock
6726 first. However, for performance reasons, implementations often favor writers over readers to
6727 avoid potential writer starvation.

6728 A read-write lock object is an implementation-defined opaque object of type
6729 **pthread_rwlock_t** as defined in `<pthread.h>`. There are two different sorts of locks
6730 associated with a read-write lock: a read lock and a write lock.

6731 The `pthread_rwlockattr_init()` function initializes a read-write lock attributes object with the
6732 default value for all the attributes defined in the implementation. After a read-write lock
6733 attributes object has been used to initialize one or more read-write locks, changes to the
6734 read-write lock attributes object, including destruction, do not affect previously initialized
6735 read-write locks.

6736 Implementations must provide at least the read-write lock attribute `process-shared`. This
6737 attribute can have the following values:

6738 **PTHREAD_PROCESS_SHARED**
6739 Any thread of any process that has access to the memory where the read-write lock
6740 resides can manipulate the read-write lock.

6741 **PTHREAD_PROCESS_PRIVATE**
6742 Only threads created within the same process as the thread that initialized the read-
6743 write lock can manipulate the read-write lock. This is the default value.

6744 The `pthread_rwlockattr_setpshared()` function is used to set the `process-shared` attribute of an
6745 initialized read-write lock attributes object while the function `pthread_rwlockattr_getpshared()`
6746 obtains the current value of the `process-shared` attribute.

6747 A read-write lock attributes object is destroyed using the `pthread_rwlockattr_destroy()`
6748 function. The effect of subsequent use of the read-write lock attributes object is undefined.

6749 A thread creates a read-write lock using the `pthread_rwlock_init()` function. The attributes of
6750 the read-write lock can be specified by the application developer; otherwise, the default
6751 implementation-defined read-write lock attributes are used if the pointer to the read-write
6752 lock attributes object is NULL. In cases where the default attributes are appropriate, the
6753 **PTHREAD_RWLOCK_INITIALIZER** macro can be used to initialize statically allocated
6754 read-write locks.

6755 A thread which wants to apply a read lock to the read-write lock can use either
6756 `pthread_rwlock_rdlock()` or `pthread_rwlock_tryrdlock()`. If `pthread_rwlock_rdlock()` is used, the
6757 thread acquires a read lock if a writer does not hold the write lock and there are no writers
6758 blocked on the write lock. If a read lock is not acquired, the calling thread blocks until it can
6759 acquire a lock. However, if `pthread_rwlock_tryrdlock()` is used, the function returns
6760 immediately with the error [EBUSY] if any thread holds a write lock or there are blocked
6761 writers waiting for the write lock.

6762 A thread which wants to apply a write lock to the read-write lock can use either of two
6763 functions: `pthread_rwlock_wrlock()` or `pthread_rwlock_trywrlock()`. If `pthread_rwlock_wrlock()`
6764 is used, the thread acquires the write lock if no other reader or writer threads hold the read-
6765 write lock. If the write lock is not acquired, the thread blocks until it can acquire the write
6766 lock. However, if `pthread_rwlock_trywrlock()` is used, the function returns immediately with

6767 the error [EBUSY] if any thread is holding either a read or a write lock.

6768 The *pthread_rwlock_unlock()* function is used to unlock a read-write lock object held by the
6769 calling thread. Results are undefined if the read-write lock is not held by the calling thread. If
6770 there are other read locks currently held on the read-write lock object, the read-write lock
6771 object remains in the read locked state but without the current thread as one of its owners. If
6772 this function releases the last read lock for this read-write lock object, the read-write lock
6773 object is put in the unlocked read state. If this function is called to release a write lock for this
6774 read-write lock object, the read-write lock object is put in the unlocked state.

6775 • Thread Concurrency Level

6776 On threads implementations that multiplex user threads onto a smaller set of kernel
6777 execution entities, the system attempts to create a reasonable number of kernel execution
6778 entities for the application upon application startup.

6779 On some implementations, these kernel entities are retained by user threads that block in the
6780 kernel. Other implementations do not *timeslice* user threads so that multiple compute-bound
6781 user threads can share a kernel thread. On such implementations, some applications may use
6782 up all the available kernel execution entities before their user-space threads are used up. The
6783 process may be left with user threads capable of doing work for the application but with no
6784 way to schedule them.

6785 The *pthread_setconcurrency()* function enables an application to request more kernel entities;
6786 that is, specify a desired concurrency level. However, this function merely provides a hint to
6787 the implementation. The implementation is free to ignore this request or to provide some
6788 other number of kernel entities. If an implementation does not multiplex user threads onto a
6789 smaller number of kernel execution entities, the *pthread_setconcurrency()* function has no
6790 effect.

6791 The *pthread_setconcurrency()* function may also have an effect on implementations where the
6792 kernel mode and user mode schedulers cooperate to ensure that ready user threads are not
6793 prevented from running by other threads blocked in the kernel.

6794 The *pthread_getconcurrency()* function always returns the value set by a previous call to
6795 *pthread_setconcurrency()*. However, if *pthread_setconcurrency()* was not previously called, this
6796 function returns zero to indicate that the threads implementation is maintaining the
6797 concurrency level.

6798 • Thread Stack Guard Size

6799 DCE threads introduced the concept of a “thread stack guard size”. Most thread
6800 implementations add a region of protected memory to a thread’s stack, commonly known as
6801 a “guard region”, as a safety measure to prevent stack pointer overflow in one thread from
6802 corrupting the contents of another thread’s stack. The default size of the guard regions
6803 attribute is {PAGESIZE} bytes and is implementation-defined.

6804 Some application developers may wish to change the stack guard size. When an application
6805 creates a large number of threads, the extra page allocated for each stack may strain system
6806 resources. In addition to the extra page of memory, the kernel’s memory manager has to keep
6807 track of the different protections on adjoining pages. When this is a problem, the application
6808 developer may request a guard size of 0 bytes to conserve system resources by eliminating
6809 stack overflow protection.

6810 Conversely an application that allocates large data structures such as arrays on the stack may
6811 wish to increase the default guard size in order to detect stack overflow. If a thread allocates
6812 two pages for a data array, a single guard page provides little protection against thread stack
6813 overflows since the thread can corrupt adjoining memory beyond the guard page.

6814 The System Interfaces volume of IEEE Std 1003.1-2001 defines a new attribute of a thread
6815 attributes object; that is, the *guardsize* attribute which allows applications to specify the size
6816 of the guard region of a thread's stack.

6817 Two functions are provided for manipulating a thread's stack guard size. The
6818 *pthread_attr_setguardsize()* function sets the thread *guardsize* attribute, and the
6819 *pthread_attr_getguardsize()* function retrieves the current value.

6820 An implementation may round up the requested guard size to a multiple of the configurable
6821 system variable {PAGESIZE}. In this case, *pthread_attr_getguardsize()* returns the guard size
6822 specified by the previous *pthread_attr_setguardsize()* function call and not the rounded up
6823 value.

6824 If an application is managing its own thread stacks using the *stackaddr* attribute, the *guardsize*
6825 attribute is ignored and no stack overflow protection is provided. In this case, it is the
6826 responsibility of the application to manage stack overflow along with stack allocation.

6827 • Parallel I/O

6828 Suppose two or more threads independently issue read requests on the same file. To read
6829 specific data from a file, a thread must first call *lseek()* to seek to the proper offset in the file,
6830 and then call *read()* to retrieve the required data. If more than one thread does this at the
6831 same time, the first thread may complete its seek call, but before it gets a chance to issue its
6832 read call a second thread may complete its seek call, resulting in the first thread accessing
6833 incorrect data when it issues its read call. One workaround is to lock the file descriptor while
6834 seeking and reading or writing, but this reduces parallelism and adds overhead.

6835 Instead, the System Interfaces volume of IEEE Std 1003.1-2001 provides two functions to
6836 make seek/read and seek/write operations atomic. The file descriptor's current offset is
6837 unchanged, thus allowing multiple read and write operations to proceed in parallel. This
6838 improves the I/O performance of threaded applications. The *pread()* function is used to do
6839 an atomic read of data from a file into a buffer. Conversely, the *pwrite()* function does an
6840 atomic write of data from a buffer to a file.

6841 B.2.9.1 Thread-Safety

6842 All functions required by IEEE Std 1003.1-2001 need to be thread-safe. Implementations have to
6843 provide internal synchronization when necessary in order to achieve this goal. In certain
6844 cases—for example, most floating-point implementations—context switch code may have to
6845 manage the writable shared state.

6846 While a read from a pipe of {PIPE_MAX}*2 bytes may not generate a single atomic and thread-
6847 safe stream of bytes, it should generate “several” (individually atomic) thread-safe streams of
6848 bytes. Similarly, while reading from a terminal device may not generate a single atomic and
6849 thread-safe stream of bytes, it should generate some finite number of (individually atomic) and
6850 thread-safe streams of bytes. That is, concurrent calls to read for a pipe, FIFO, or terminal device
6851 are not allowed to result in corrupting the stream of bytes or other internal data. However,
6852 *read()*, in these cases, is not required to return a single contiguous and atomic stream of bytes.

6853 It is not required that all functions provided by IEEE Std 1003.1-2001 be either async-cancel-safe
6854 or async-signal-safe.

6855 As it turns out, some functions are inherently not thread-safe; that is, their interface
6856 specifications preclude reentrancy. For example, some functions (such as *asctime()*) return a
6857 pointer to a result stored in memory space allocated by the function on a per-process basis. Such
6858 a function is not thread-safe, because its result can be overwritten by successive invocations.
6859 Other functions, while not inherently non-thread-safe, may be implemented in ways that lead to

6860 them not being thread-safe. For example, some functions (such as *rand()*) store state information
6861 (such as a seed value, which survives multiple function invocations) in memory space allocated
6862 by the function on a per-process basis. The implementation of such a function is not thread-safe
6863 if the implementation fails to synchronize invocations of the function and thus fails to protect
6864 the state information. The problem is that when the state information is not protected,
6865 concurrent invocations can interfere with one another (for example, applications using *rand()*
6866 may see the same seed value).

6867 *Thread-Safety and Locking of Existing Functions*

6868 Originally, POSIX.1 was not designed to work in a multi-threaded environment, and some
6869 implementations of some existing functions will not work properly when executed concurrently.
6870 To provide routines that will work correctly in an environment with threads (“thread-safe”), two
6871 problems need to be solved:

- 6872 1. Routines that maintain or return pointers to static areas internal to the routine (which may
6873 now be shared) need to be modified. The routines *ttyname()* and *localtime()* are examples.
- 6874 2. Routines that access data space shared by more than one thread need to be modified. The
6875 *malloc()* function and the *stdio* family routines are examples.

6876 There are a variety of constraints on these changes. The first is compatibility with the existing
6877 versions of these functions—non-thread-safe functions will continue to be in use for some time,
6878 as the original interfaces are used by existing code. Another is that the new thread-safe versions
6879 of these functions represent as small a change as possible over the familiar interfaces provided
6880 by the existing non-thread-safe versions. The new interfaces should be independent of any
6881 particular threads implementation. In particular, they should be thread-safe without depending
6882 on explicit thread-specific memory. Finally, there should be minimal performance penalty due to
6883 the changes made to the functions.

6884 It is intended that the list of functions from POSIX.1 that cannot be made thread-safe and for
6885 which corrected versions are provided be complete.

6886 *Thread-Safety and Locking Solutions*

6887 Many of the POSIX.1 functions were thread-safe and did not change at all. However, some
6888 functions (for example, the math functions typically found in **libm**) are not thread-safe because
6889 of writable shared global state. For instance, in IEEE Std 754-1985 floating-point
6890 implementations, the computation modes and flags are global and shared.

6891 Some functions are not thread-safe because a particular implementation is not reentrant,
6892 typically because of a non-essential use of static storage. These require only a new
6893 implementation.

6894 Thread-safe libraries are useful in a wide range of parallel (and asynchronous) programming
6895 environments, not just within pthreads. In order to be used outside the context of pthreads,
6896 however, such libraries still have to use some synchronization method. These could either be
6897 independent of the pthread synchronization operations, or they could be a subset of the pthread
6898 interfaces. Either method results in thread-safe library implementations that can be used without
6899 the rest of pthreads.

6900 Some functions, such as the *stdio* family interface and dynamic memory allocation functions
6901 such as *malloc()*, are inter-dependent routines that share resources (for example, buffers) across
6902 related calls. These require synchronization to work correctly, but they do not require any
6903 change to their external (user-visible) interfaces.

6904 In some cases, such as *getc()* and *putc()*, adding synchronization is likely to create an
6905 unacceptable performance impact. In this case, slower thread-safe synchronized functions are to

6906 be provided, but the original, faster (but unsafe) functions (which may be implemented as
 6907 macros) are retained under new names. Some additional special-purpose synchronization
 6908 facilities are necessary for these macros to be usable in multi-threaded programs. This also
 6909 requires changes in `<stdio.h>`.

6910 The other common reason that functions are unsafe is that they return a pointer to static storage,
 6911 making the functions non-thread-safe. This has to be changed, and there are three natural
 6912 choices:

6913 1. Return a pointer to thread-specific storage

6914 This could incur a severe performance penalty on those architectures with a costly
 6915 implementation of the thread-specific data interface.

6916 A variation on this technique is to use `malloc()` to allocate storage for the function output
 6917 and return a pointer to this storage. This technique may also have an undesirable
 6918 performance impact, however, and a simplistic implementation requires that the user
 6919 program explicitly free the storage object when it is no longer needed. This technique is
 6920 used by some existing POSIX.1 functions. With careful implementation for infrequently
 6921 used functions, there may be little or no performance or storage penalty, and the
 6922 maintenance of already-standardized interfaces is a significant benefit.

6923 2. Return the actual value computed by the function

6924 This technique can only be used with functions that return pointers to structures—routines
 6925 that return character strings would have to wrap their output in an enclosing structure in
 6926 order to return the output on the stack. There is also a negative performance impact
 6927 inherent in this solution in that the output value has to be copied twice before it can be
 6928 used by the calling function: once from the called routine's local buffers to the top of the
 6929 stack, then from the top of the stack to the assignment target. Finally, many older
 6930 compilers cannot support this technique due to a historical tendency to use internal static
 6931 buffers to deliver the results of structure-valued functions.

6932 3. Have the caller pass the address of a buffer to contain the computed value

6933 The only disadvantage of this approach is that extra arguments have to be provided by the
 6934 calling program. It represents the most efficient solution to the problem, however, and,
 6935 unlike the `malloc()` technique, it is semantically clear.

6936 There are some routines (often groups of related routines) whose interfaces are inherently non-
 6937 thread-safe because they communicate across multiple function invocations by means of static
 6938 memory locations. The solution is to redesign the calls so that they are thread-safe, typically by
 6939 passing the needed data as extra parameters. Unfortunately, this may require major changes to
 6940 the interface as well.

6941 A floating-point implementation using IEEE Std 754-1985 is a case in point. A less problematic
 6942 example is the `rand48` family of pseudo-random number generators. The functions `getgrgid()`,
 6943 `getgrnam()`, `getpwnam()`, and `getpwuid()` are another such case.

6944 The problems with `errno` are discussed in **Alternative Solutions for Per-Thread `errno`** (on page
 6945 94).

6946 Some functions can be thread-safe or not, depending on their arguments. These include the
 6947 `tmpnam()` and `ctermid()` functions. These functions have pointers to character strings as
 6948 arguments. If the pointers are not NULL, the functions store their results in the character string;
 6949 however, if the pointers are NULL, the functions store their results in an area that may be static
 6950 and thus subject to overwriting by successive calls. These should only be called by multi-thread
 6951 applications when their arguments are non-NULL.

6952 *Asynchronous Safety and Thread-Safety*

6953 A floating-point implementation has many modes that effect rounding and other aspects of
6954 computation. Functions in some math library implementations may change the computation
6955 modes for the duration of a function call. If such a function call is interrupted by a signal or
6956 cancelation, the floating-point state is not required to be protected.

6957 There is a significant cost to make floating-point operations async-cancel-safe or async-signal-
6958 safe; accordingly, neither form of async safety is required.

6959 *Functions Returning Pointers to Static Storage*

6960 For those functions that are not thread-safe because they return values in fixed size statically
6961 allocated structures, alternate “_r” forms are provided that pass a pointer to an explicit result
6962 structure. Those that return pointers into library-allocated buffers have forms provided with
6963 explicit buffer and length parameters.

6964 For functions that return pointers to library-allocated buffers, it makes sense to provide “_r”
6965 versions that allow the application control over allocation of the storage in which results are
6966 returned. This allows the state used by these functions to be managed on an application-specific
6967 basis, supporting per-thread, per-process, or other application-specific sharing relationships.

6968 Early proposals had provided “_r” versions for functions that returned pointers to variable-size
6969 buffers without providing a means for determining the required buffer size. This would have
6970 made using such functions exceedingly clumsy, potentially requiring iteratively calling them
6971 with increasingly larger guesses for the amount of storage required. Hence, *sysconf()* variables
6972 have been provided for such functions that return the maximum required buffer size.

6973 Thus, the rule that has been followed by IEEE Std 1003.1-2001 when adapting single-threaded
6974 non-thread-safe functions is as follows: all functions returning pointers to library-allocated
6975 storage should have “_r” versions provided, allowing the application control over the storage
6976 allocation. Those with variable-sized return values accept both a buffer address and a length
6977 parameter. The *sysconf()* variables are provided to supply the appropriate buffer sizes when
6978 required. Implementors are encouraged to apply the same rule when adapting their own existing
6979 functions to a pthreads environment.

6980 *B.2.9.2 Thread IDs*

6981 Separate applications should communicate through well-defined interfaces and should not
6982 depend on each other's implementation. For example, if a programmer decides to rewrite the *sort*
6983 utility using multiple threads, it should be easy to do this so that the interface to the *sort*
6984 utility does not change. Consider that if the user causes SIGINT to be generated while the *sort*
6985 utility is running, keeping the same interface means that the entire *sort* utility is killed, not just one of its
6986 threads. As another example, consider a realtime application that manages a reactor. Such an
6987 application may wish to allow other applications to control the priority at which it watches the
6988 control rods. One technique to accomplish this is to write the ID of the thread watching the
6989 control rods into a file and allow other programs to change the priority of that thread as they see
6990 fit. A simpler technique is to have the reactor process accept IPCs (Interprocess Communication
6991 messages) from other processes, telling it at a semantic level what priority the program should
6992 assign to watching the control rods. This allows the programmer greater flexibility in the
6993 implementation. For example, the programmer can change the implementation from having one
6994 thread per rod to having one thread watching all of the rods without changing the interface.
6995 Having threads live inside the process means that the implementation of a process is invisible to
6996 outside processes (excepting debuggers and system management tools).

6997 Threads do not provide a protection boundary. Every thread model allows threads to share
6998 memory with other threads and encourages this sharing to be widespread. This means that one

6999 thread can wipe out memory that is needed for the correct functioning of other threads that are
7000 sharing its memory. Consequently, providing each thread with its own user and/or group IDs
7001 would not provide a protection boundary between threads sharing memory.

7002 *B.2.9.3 Thread Mutexes*

7003 There is no additional rationale provided for this section.

7004 *B.2.9.4 Thread Scheduling*

7005 • Scheduling Implementation Models

7006 The following scheduling implementation models are presented in terms of threads and
7007 “kernel entities”. This is to simplify exposition of the models, and it does not imply that an
7008 implementation actually has an identifiable “kernel entity”.

7009 A kernel entity is not defined beyond the fact that it has scheduling attributes that are used to
7010 resolve contention with other kernel entities for execution resources. A kernel entity may be
7011 thought of as an envelope that holds a thread or a separate kernel thread. It is not a
7012 conventional process, although it shares with the process the attribute that it has a single
7013 thread of control; it does not necessarily imply an address space, open files, and so on. It is
7014 better thought of as a primitive facility upon which conventional processes and threads may
7015 be constructed.

7016 — System Thread Scheduling Model

7017 This model consists of one thread per kernel entity. The kernel entity is solely responsible
7018 for scheduling thread execution on one or more processors. This model schedules all
7019 threads against all other threads in the system using the scheduling attributes of the
7020 thread.

7021 — Process Scheduling Model

7022 A generalized process scheduling model consists of two levels of scheduling. A threads
7023 library creates a pool of kernel entities, as required, and schedules threads to run on them
7024 using the scheduling attributes of the threads. Typically, the size of the pool is a function
7025 of the simultaneously runnable threads, not the total number of threads. The kernel then
7026 schedules the kernel entities onto processors according to their scheduling attributes,
7027 which are managed by the threads library. This set model potentially allows a wide range
7028 of mappings between threads and kernel entities.

7029 • System and Process Scheduling Model Performance

7030 There are a number of important implications on the performance of applications using these
7031 scheduling models. The process scheduling model potentially provides lower overhead for
7032 making scheduling decisions, since there is no need to access kernel-level information or
7033 functions and the set of schedulable entities is smaller (only the threads within the process).

7034 On the other hand, since the kernel is also making scheduling decisions regarding the system
7035 resources under its control (for example, CPU(s), I/O devices, memory), decisions that do
7036 not take thread scheduling parameters into account can result in unspecified delays for
7037 realtime application threads, causing them to miss maximum response time limits.

7038 • Rate Monotonic Scheduling

7039 Rate monotonic scheduling was considered, but rejected for standardization in the context of
7040 pthreads. A sporadic server policy is included.

7041 • Scheduling Options

7042 In IEEE Std 1003.1-2001, the basic thread scheduling functions are defined under the Threads
7043 option, so that they are required of all threads implementations. However, there are no
7044 specific scheduling policies required by this option to allow for conforming thread
7045 implementations that are not targeted to realtime applications.

7046 Specific standard scheduling policies are defined to be under the Thread Execution
7047 Scheduling option, and they are specifically designed to support realtime applications by
7048 providing predictable resource-sharing sequences. The name of this option was chosen to
7049 emphasize that this functionality is defined as appropriate for realtime applications that
7050 require simple priority-based scheduling.

7051 It is recognized that these policies are not necessarily satisfactory for some multi-processor
7052 implementations, and work is ongoing to address a wider range of scheduling behaviors. The
7053 interfaces have been chosen to create abundant opportunity for future scheduling policies to
7054 be implemented and standardized based on this interface. In order to standardize a new
7055 scheduling policy, all that is required (from the standpoint of thread scheduling attributes) is
7056 to define a new policy name, new members of the thread attributes object, and functions to
7057 set these members when the scheduling policy is equal to the new value.

7058 **Scheduling Contention Scope**

7059 In order to accommodate the requirement for realtime response, each thread has a scheduling
7060 contention scope attribute. Threads with a system scheduling contention scope have to be
7061 scheduled with respect to all other threads in the system. These threads are usually bound to a
7062 single kernel entity that reflects their scheduling attributes and are directly scheduled by the
7063 kernel.

7064 Threads with a process scheduling contention scope need be scheduled only with respect to the
7065 other threads in the process. These threads may be scheduled within the process onto a pool of
7066 kernel entities. The implementation is also free to bind these threads directly to kernel entities
7067 and let them be scheduled by the kernel. Process scheduling contention scope allows the
7068 implementation the most flexibility and is the default if both contention scopes are supported
7069 and none is specified.

7070 Thus, the choice by implementors to provide one or the other (or both) of these scheduling
7071 models is driven by the need of their supported application domains for worst-case (that is,
7072 realtime) response, or average-case (non-realtime) response.

7073 **Scheduling Allocation Domain**

7074 The SCHED_FIFO and SCHED_RR scheduling policies take on different characteristics on a
7075 multi-processor. Other scheduling policies are also subject to changed behavior when executed
7076 on a multi-processor. The concept of scheduling allocation domain determines the set of
7077 processors on which the threads of an application may run. By considering the application's
7078 processor scheduling allocation domain for its threads, scheduling policies can be defined in
7079 terms of their behavior for varying processor scheduling allocation domain values. It is
7080 conceivable that not all scheduling allocation domain sizes make sense for all scheduling
7081 policies on all implementations. The concept of scheduling allocation domain, however, is a
7082 useful tool for the description of multi-processor scheduling policies.

7083 The “process control” approach to scheduling obtains significant performance advantages from
7084 dynamic scheduling allocation domain sizes when it is applicable.

7085 Non-Uniform Memory Access (NUMA) multi-processors may use a system scheduling structure
7086 that involves reassignment of threads among scheduling allocation domains. In NUMA

7087 machines, a natural model of scheduling is to match scheduling allocation domains to clusters of
7088 processors. Load balancing in such an environment requires changing the scheduling allocation
7089 domain to which a thread is assigned.

7090 **Scheduling Documentation**

7091 Implementation-provided scheduling policies need to be completely documented in order to be
7092 useful. This documentation includes a description of the attributes required for the policy, the
7093 scheduling interaction of threads running under this policy and all other supported policies, and
7094 the effects of all possible values for processor scheduling allocation domain. Note that for the
7095 implementor wishing to be minimally-compliant, it is (minimally) acceptable to define the
7096 behavior as undefined.

7097 **Scheduling Contention Scope Attribute**

7098 The scheduling contention scope defines how threads compete for resources. Within
7099 IEEE Std 1003.1-2001, scheduling contention scope is used to describe only how threads are
7100 scheduled in relation to one another in the system. That is, either they are scheduled against all
7101 other threads in the system (“system scope”) or only against those threads in the process
7102 (“process scope”). In fact, scheduling contention scope may apply to additional resources,
7103 including virtual timers and profiling, which are not currently considered by
7104 IEEE Std 1003.1-2001.

7105 **Mixed Scopes**

7106 If only one scheduling contention scope is supported, the scheduling decision is straightforward.
7107 To perform the processor scheduling decision in a mixed scope environment, it is necessary to
7108 map the scheduling attributes of the thread with process-wide contention scope to the same
7109 attribute space as the thread with system-wide contention scope.

7110 Since a conforming implementation has to support one and may support both scopes, it is useful
7111 to discuss the effects of such choices with respect to example applications. If an implementation
7112 supports both scopes, mixing scopes provides a means of better managing system-level (that is,
7113 kernel-level) and library-level resources. In general, threads with system scope will require the
7114 resources of a separate kernel entity in order to guarantee the scheduling semantics. On the
7115 other hand, threads with process scope can share the resources of a kernel entity while
7116 maintaining the scheduling semantics.

7117 The application is free to create threads with dedicated kernel resources, and other threads that
7118 multiplex kernel resources. Consider the example of a window server. The server allocates two
7119 threads per widget: one thread manages the widget user interface (including drawing), while the
7120 other thread takes any required application action. This allows the widget to be “active” while
7121 the application is computing. A screen image may be built from thousands of widgets. If each of
7122 these threads had been created with system scope, then most of the kernel-level resources might
7123 be wasted, since only a few widgets are active at any one time. In addition, mixed scope is
7124 particularly useful in a window server where one thread with high priority and system scope
7125 handles the mouse so that it tracks well. As another example, consider a database server. For
7126 each of the hundreds or thousands of clients supported by a large server, an equivalent number
7127 of threads will have to be created. If each of these threads were system scope, the consequences
7128 would be the same as for the window server example above. However, the server could be
7129 constructed so that actual retrieval of data is done by several dedicated threads. Dedicated
7130 threads that do work for all clients frequently justify the added expense of system scope. If it
7131 were not permissible to mix system and process threads in the same process, this type of
7132 solution would not be possible.

7133 Dynamic Thread Scheduling Parameters Access

7134 In many time-constrained applications, there is no need to change the scheduling attributes
7135 dynamically during thread or process execution, since the general use of these attributes is to
7136 reflect directly the time constraints of the application. Since these time constraints are generally
7137 imposed to meet higher-level system requirements, such as accuracy or availability, they
7138 frequently should remain unchanged during application execution.

7139 However, there are important situations in which the scheduling attributes should be changed.
7140 Generally, this will occur when external environmental conditions exist in which the time
7141 constraints change. Consider, for example, a space vehicle major mode change, such as the
7142 change from ascent to descent mode, or the change from the space environment to the
7143 atmospheric environment. In such cases, the frequency with which many of the sensors or
7144 actuators need to be read or written will change, which will necessitate a priority change. In
7145 other cases, even the existence of a time constraint might be temporary, necessitating not just a
7146 priority change, but also a policy change for ongoing threads or processes. For this reason, it is
7147 critical that the interface should provide functions to change the scheduling parameters
7148 dynamically, but, as with many of the other realtime functions, it is important that applications
7149 use them properly to avoid the possibility of unnecessarily degrading performance.

7150 In providing functions for dynamically changing the scheduling behavior of threads, there were
7151 two options: provide functions to get and set the individual scheduling parameters of threads, or
7152 provide a single interface to get and set all the scheduling parameters for a given thread
7153 simultaneously. Both approaches have merit. Access functions for individual parameters allow
7154 simpler control of thread scheduling for simple thread scheduling parameters. However, a single
7155 function for setting all the parameters for a given scheduling policy is required when first setting
7156 that scheduling policy. Since the single all-encompassing functions are required, it was decided
7157 to leave the interface as minimal as possible. Note that simpler functions (such as
7158 *pthread_setprio()* for threads running under the priority-based schedulers) can be easily defined
7159 in terms of the all-encompassing functions.

7160 If the *pthread_setschedparam()* function executes successfully, it will have set all of the scheduling
7161 parameter values indicated in *param*; otherwise, none of the scheduling parameters will have
7162 been modified. This is necessary to ensure that the scheduling of this and all other threads
7163 continues to be consistent in the presence of an erroneous scheduling parameter.

7164 The [EPERM] error value is included in the list of possible *pthread_setschedparam()* error returns
7165 as a reflection of the fact that the ability to change scheduling parameters increases risks to the
7166 implementation and application performance if the scheduling parameters are changed
7167 improperly. For this reason, and based on some existing practice, it was felt that some
7168 implementations would probably choose to define specific permissions for changing either a
7169 thread's own or another thread's scheduling parameters. IEEE Std 1003.1-2001 does not include
7170 portable methods for setting or retrieving permissions, so any such use of permissions is
7171 completely unspecified.

7172 Mutex Initialization Scheduling Attributes

7173 In a priority-driven environment, a direct use of traditional primitives like mutexes and
7174 condition variables can lead to unbounded priority inversion, where a higher priority thread can
7175 be blocked by a lower priority thread, or set of threads, for an unbounded duration of time. As a
7176 result, it becomes impossible to guarantee thread deadlines. Priority inversion can be bounded
7177 and minimized by the use of priority inheritance protocols. This allows thread deadlines to be
7178 guaranteed even in the presence of synchronization requirements.

7179 Two useful but simple members of the family of priority inheritance protocols are the basic
7180 priority inheritance protocol and the priority ceiling protocol emulation. Under the Basic Priority

7181 Inheritance protocol (governed by the Thread Priority Inheritance option), a thread that is
7182 blocking higher priority threads executes at the priority of the highest priority thread that it
7183 blocks. This simple mechanism allows priority inversion to be bounded by the duration of
7184 critical sections and makes timing analysis possible.

7185 Under the Priority Ceiling Protocol Emulation protocol (governed by the Thread Priority
7186 Protection option), each mutex has a priority ceiling, usually defined as the priority of the
7187 highest priority thread that can lock the mutex. When a thread is executing inside critical
7188 sections, its priority is unconditionally increased to the highest of the priority ceilings of all the
7189 mutexes owned by the thread. This protocol has two very desirable properties in uni-processor
7190 systems. First, a thread can be blocked by a lower priority thread for at most the duration of one
7191 single critical section. Furthermore, when the protocol is correctly used in a single processor, and
7192 if threads do not become blocked while owning mutexes, mutual deadlocks are prevented.

7193 The priority ceiling emulation can be extended to multiple processor environments, in which
7194 case the values of the priority ceilings will be assigned depending on the kind of mutex that is
7195 being used: local to only one processor, or global, shared by several processors. Local priority
7196 ceilings will be assigned the usual way, equal to the priority of the highest priority thread that
7197 may lock that mutex. Global priority ceilings will usually be assigned a priority level higher than
7198 all the priorities assigned to any of the threads that reside in the involved processors to avoid the
7199 effect called remote blocking.

7200 **Change the Priority Ceiling of a Mutex**

7201 In order for the priority protect protocol to exhibit its desired properties of bounding priority
7202 inversion and avoidance of deadlock, it is critical that the ceiling priority of a mutex be the same
7203 as the priority of the highest thread that can ever hold it, or higher. Thus, if the priorities of the
7204 threads using such mutexes never change dynamically, there is no need ever to change the
7205 priority ceiling of a mutex.

7206 However, if a major system mode change results in an altered response time requirement for one
7207 or more application threads, their priority has to change to reflect it. It will occasionally be the
7208 case that the priority ceilings of mutexes held also need to change. While changing priority
7209 ceilings should generally be avoided, it is important that IEEE Std 1003.1-2001 provide these
7210 interfaces for those cases in which it is necessary.

7211 *B.2.9.5 Thread Cancellation*

7212 Many existing threads packages have facilities for canceling an operation or canceling a thread.
7213 These facilities are used for implementing user requests (such as the CANCEL button in a
7214 window-based application), for implementing OR parallelism (for example, telling the other
7215 threads to stop working once one thread has found a forced mate in a parallel chess program), or
7216 for implementing the ABORT mechanism in Ada.

7217 POSIX programs traditionally have used the signal mechanism combined with either *longjmp()*
7218 or polling to cancel operations. Many POSIX programmers have trouble using these facilities to
7219 solve their problems efficiently in a single-threaded process. With the introduction of threads,
7220 these solutions become even more difficult to use.

7221 The main issues with implementing a cancellation facility are specifying the operation to be
7222 canceled, cleanly releasing any resources allocated to that operation, controlling when the target
7223 notices that it has been canceled, and defining the interaction between asynchronous signals and
7224 cancellation.

7225 Specifying the Operation to Cancel

7226 Consider a thread that calls through five distinct levels of program abstraction and then, inside
7227 the lowest-level abstraction, calls a function that suspends the thread. (An abstraction boundary
7228 is a layer at which the client of the abstraction sees only the service being provided and can
7229 remain ignorant of the implementation. Abstractions are often layered, each level of abstraction
7230 being a client of the lower-level abstraction and implementing a higher-level abstraction.)
7231 Depending on the semantics of each abstraction, one could imagine wanting to cancel only the
7232 call that causes suspension, only the bottom two levels, or the operation being done by the entire
7233 thread. Canceling operations at a finer grain than the entire thread is difficult because threads
7234 are active and they may be run in parallel on a multi-processor. By the time one thread can make
7235 a request to cancel an operation, the thread performing the operation may have completed that
7236 operation and gone on to start another operation whose cancelation is not desired. Thread IDs
7237 are not reused until the thread has exited, and either it was created with the *Attr detachstate*
7238 attribute set to *PTHREAD_CREATE_DETACHED* or the *pthread_join()* or *pthread_detach()*
7239 function has been called for that thread. Consequently, a thread cancelation will never be
7240 misdirected when the thread terminates. For these reasons, the canceling of operations is done at
7241 the granularity of the thread. Threads are designed to be inexpensive enough so that a separate
7242 thread may be created to perform each separately cancelable operation; for example, each
7243 possibly long running user request.

7244 For cancelation to be used in existing code, cancelation scopes and handlers will have to be
7245 established for code that needs to release resources upon cancelation, so that it follows the
7246 programming discipline described in the text.

7247 A Special Signal Versus a Special Interface

7248 Two different mechanisms were considered for providing the cancelation interfaces. The first
7249 was to provide an interface to direct signals at a thread and then to define a special signal that
7250 had the required semantics. The other alternative was to use a special interface that delivered the
7251 correct semantics to the target thread.

7252 The solution using signals produced a number of problems. It required the implementation to
7253 provide cancelation in terms of signals whereas a perfectly valid (and possibly more efficient)
7254 implementation could have both layered on a low-level set of primitives. There were so many
7255 exceptions to the special signal (it cannot be used with *kill()*, no POSIX.1 interfaces can be used
7256 with it) that it was clearly not a valid signal. Its semantics on delivery were also completely
7257 different from any existing POSIX.1 signal. As such, a special interface that did not mandate the
7258 implementation and did not confuse the semantics of signals and cancelation was felt to be the
7259 better solution.

7260 Races Between Cancelation and Resuming Execution

7261 Due to the nature of cancelation, there is generally no synchronization between the thread
7262 requesting the cancelation of a blocked thread and events that may cause that thread to resume
7263 execution. For this reason, and because excess serialization hurts performance, when both an
7264 event that a thread is waiting for has occurred and a cancelation request has been made and
7265 cancelation is enabled, IEEE Std 1003.1-2001 explicitly allows the implementation to choose
7266 between returning from the blocking call or acting on the cancelation request.

7267 **Interaction of Cancellation with Asynchronous Signals**

7268 A typical use of cancellation is to acquire a lock on some resource and to establish a cancellation
7269 cleanup handler for releasing the resource when and if the thread is canceled.

7270 A correct and complete implementation of cancellation in the presence of asynchronous signals
7271 requires considerable care. An implementation has to push a cancellation cleanup handler on the
7272 cancellation cleanup stack while maintaining the integrity of the stack data structure. If an
7273 asynchronously-generated signal is posted to the thread during a stack operation, the signal
7274 handler cannot manipulate the cancellation cleanup stack. As a consequence, asynchronous
7275 signal handlers may not cancel threads or otherwise manipulate the cancellation state of a thread.
7276 Threads may, of course, be canceled by another thread that used a *sigwait()* function to wait
7277 synchronously for an asynchronous signal.

7278 In order for cancellation to function correctly, it is required that asynchronous signal handlers not
7279 change the cancellation state. This requires that some elements of existing practice, such as using
7280 *longjmp()* to exit from an asynchronous signal handler implicitly, be prohibited in cases where
7281 the integrity of the cancellation state of the interrupt thread cannot be ensured.

7282 **Thread Cancellation Overview**

7283 • Cancellability States

7284 The three possible cancellability states (disabled, deferred, and asynchronous) are encoded
7285 into two separate bits ((disable, enable) and (deferred, asynchronous)) to allow them to be
7286 changed and restored independently. For instance, short code sequences that will not block
7287 sometimes disable cancellability on entry and restore the previous state upon exit. Likewise,
7288 long or unbounded code sequences containing no convenient explicit cancellation points will
7289 sometimes set the cancellability type to asynchronous on entry and restore the previous value
7290 upon exit.

7291 • Cancellation Points

7292 Cancellation points are points inside of certain functions where a thread has to act on any
7293 pending cancellation request when cancellability is enabled, if the function would block. As
7294 with checking for signals, operations need only check for pending cancellation requests when
7295 the operation is about to block indefinitely.

7296 The idea was considered of allowing implementations to define whether blocking calls such
7297 as *read()* should be cancellation points. It was decided that it would adversely affect the
7298 design of conforming applications if blocking calls were not cancellation points because
7299 threads could be left blocked in an uncancelable state.

7300 There are several important blocking routines that are specifically not made cancellation
7301 points:

7302 — *pthread_mutex_lock()*

7303 If *pthread_mutex_lock()* were a cancellation point, every routine that called it would also
7304 become a cancellation point (that is, any routine that touched shared state would
7305 automatically become a cancellation point). For example, *malloc()*, *free()*, and *rand()*
7306 would become cancellation points under this scheme. Having too many cancellation points
7307 makes programming very difficult, leading to either much disabling and restoring of
7308 cancellability or much difficulty in trying to arrange for reliable cleanup at every possible
7309 place.

7310 Since *pthread_mutex_lock()* is not a cancellation point, threads could result in being
7311 blocked uninterruptibly for long periods of time if mutexes were used as a general

7312 synchronization mechanism. As this is normally not acceptable, mutexes should only be
 7313 used to protect resources that are held for small fixed lengths of time where not being
 7314 able to be canceled will not be a problem. Resources that need to be held exclusively for
 7315 long periods of time should be protected with condition variables.

7316 — *pthread_barrier_wait()*

7317 Canceling a barrier wait will render a barrier unusable. Similar to a barrier timeout (which
 7318 the standard developers rejected), there is no way to guarantee the consistency of a
 7319 barrier's internal data structures if a barrier wait is canceled.

7320 — *pthread_spin_lock()*

7321 As with mutexes, spin locks should only be used to protect resources that are held for
 7322 small fixed lengths of time where not being cancelable will not be a problem.

7323 Every library routine should specify whether or not it includes any cancellation points.
 7324 Typically, only those routines that may block or compute indefinitely need to include
 7325 cancellation points.

7326 Correctly coded routines only reach cancellation points after having set up a cancellation
 7327 cleanup handler to restore invariants if the thread is canceled at that point. Being cancelable
 7328 only at specified cancellation points allows programmers to keep track of actions needed in a
 7329 cancellation cleanup handler more easily. A thread should only be made asynchronously
 7330 cancelable when it is not in the process of acquiring or releasing resources or otherwise in a
 7331 state from which it would be difficult or impossible to recover.

7332 • Thread Cancellation Cleanup Handlers

7333 The cancellation cleanup handlers provide a portable mechanism, easy to implement, for
 7334 releasing resources and restoring invariants. They are easier to use than signal handlers
 7335 because they provide a stack of cancellation cleanup handlers rather than a single handler,
 7336 and because they have an argument that can be used to pass context information to the
 7337 handler.

7338 The alternative to providing these simple cancellation cleanup handlers (whose only use is for
 7339 cleaning up when a thread is canceled) is to define a general exception package that could be
 7340 used for handling and cleaning up after hardware traps and software-detected errors. This
 7341 was too far removed from the charter of providing threads to handle asynchrony. However,
 7342 it is an explicit goal of IEEE Std 1003.1-2001 to be compatible with existing exception facilities
 7343 and languages having exceptions.

7344 The interaction of this facility and other procedure-based or language-level exception
 7345 facilities is unspecified in this version of IEEE Std 1003.1-2001. However, it is intended that it
 7346 be possible for an implementation to define the relationship between these cancellation
 7347 cleanup handlers and Ada, C++, or other language-level exception handling facilities.

7348 It was suggested that the cancellation cleanup handlers should also be called when the
 7349 process exits or calls the *exec* function. This was rejected partly due to the performance
 7350 problem caused by having to call the cancellation cleanup handlers of every thread before the
 7351 operation could continue. The other reason was that the only state expected to be cleaned up
 7352 by the cancellation cleanup handlers would be the intraprocess state. Any handlers that are to
 7353 clean up the interprocess state would be registered with *atexit()*. There is the orthogonal
 7354 problem that the *exec* functions do not honor the *atexit()* handlers, but resolving this is
 7355 beyond the scope of IEEE Std 1003.1-2001.

7356 • Async-Cancel Safety

7357 A function is said to be async-cancel-safe if it is written in such a way that entering the
 7358 function with asynchronous cancelability enabled will not cause any invariants to be
 7359 violated, even if a cancelation request is delivered at any arbitrary instruction. Functions that
 7360 are async-cancel-safe are often written in such a way that they need to acquire no resources
 7361 for their operation and the visible variables that they may write are strictly limited.

7362 Any routine that gets a resource as a side effect cannot be made async-cancel-safe (for
 7363 example, *malloc()*). If such a routine were called with asynchronous cancelability enabled, it
 7364 might acquire the resource successfully, but as it was returning to the client, it could act on a
 7365 cancelation request. In such a case, the application would have no way of knowing whether
 7366 the resource was acquired or not.

7367 Indeed, because many interesting routines cannot be made async-cancel-safe, most library
 7368 routines in general are not async-cancel-safe. Every library routine should specify whether or
 7369 not it is async-cancel safe so that programmers know which routines can be called from code
 7370 that is asynchronously cancelable.

7371 *B.2.9.6 Thread Read-Write Locks*7372 **Background**

7373 Read-write locks are often used to allow parallel access to data on multi-processors, to avoid
 7374 context switches on uni-processors when multiple threads access the same data, and to protect
 7375 data structures that are frequently accessed (that is, read) but rarely updated (that is, written).
 7376 The in-core representation of a file system directory is a good example of such a data structure.
 7377 One would like to achieve as much concurrency as possible when searching directories, but limit
 7378 concurrent access when adding or deleting files.

7379 Although read-write locks can be implemented with mutexes and condition variables, such
 7380 implementations are significantly less efficient than is possible. Therefore, this synchronization
 7381 primitive is included in IEEE Std 1003.1-2001 for the purpose of allowing more efficient
 7382 implementations in multi-processor systems.

7383 **Queuing of Waiting Threads**

7384 The *pthread_rwlock_unlock()* function description states that one writer or one or more readers
 7385 must acquire the lock if it is no longer held by any thread as a result of the call. However, the
 7386 function does not specify which thread(s) acquire the lock, unless the Thread Execution
 7387 Scheduling option is supported.

7388 The standard developers considered the issue of scheduling with respect to the queuing of
 7389 threads blocked on a read-write lock. The question turned out to be whether
 7390 IEEE Std 1003.1-2001 should require priority scheduling of read-write locks for threads whose
 7391 execution scheduling policy is priority-based (for example, SCHED_FIFO or SCHED_RR). There
 7392 are tradeoffs between priority scheduling, the amount of concurrency achievable among readers,
 7393 and the prevention of writer and/or reader starvation.

7394 For example, suppose one or more readers hold a read-write lock and the following threads
 7395 request the lock in the listed order:

7396 pthread_rwlock_wrlock() - Low priority thread writer_a
 7397 pthread_rwlock_rdlock() - High priority thread reader_a
 7398 pthread_rwlock_rdlock() - High priority thread reader_b
 7399 pthread_rwlock_rdlock() - High priority thread reader_c

7400 When the lock becomes available, should *writer_a* block the high priority readers? Or, suppose a
 7401 read-write lock becomes available and the following are queued:

7402 pthread_rwlock_rdlock() - Low priority thread reader_a
 7403 pthread_rwlock_rdlock() - Low priority thread reader_b
 7404 pthread_rwlock_rdlock() - Low priority thread reader_c
 7405 pthread_rwlock_wrlock() - Medium priority thread writer_a
 7406 pthread_rwlock_rdlock() - High priority thread reader_d

7407 If priority scheduling is applied then *reader_d* would acquire the lock and *writer_a* would block
 7408 the remaining readers. But should the remaining readers also acquire the lock to increase
 7409 concurrency? The solution adopted takes into account that when the Thread Execution
 7410 Scheduling option is supported, high priority threads may in fact starve low priority threads (the
 7411 application developer is responsible in this case for designing the system in such a way that this
 7412 starvation is avoided). Therefore, IEEE Std 1003.1-2001 specifies that high priority readers take
 7413 precedence over lower priority writers. However, to prevent writer starvation from threads of
 7414 the same or lower priority, writers take precedence over readers of the same or lower priority.

7415 Priority inheritance mechanisms are non-trivial in the context of read-write locks. When a high
 7416 priority writer is forced to wait for multiple readers, for example, it is not clear which subset of
 7417 the readers should inherit the writer's priority. Furthermore, the internal data structures that
 7418 record the inheritance must be accessible to all readers, and this implies some sort of
 7419 serialization that could negate any gain in parallelism achieved through the use of multiple
 7420 readers in the first place. Finally, existing practice does not support the use of priority
 7421 inheritance for read-write locks. Therefore, no specification of priority inheritance or priority
 7422 ceiling is attempted. If reliable priority-scheduled synchronization is absolutely required, it can
 7423 always be obtained through the use of mutexes.

7424 **Comparison to *fcntl()* Locks**

7425 The read-write locks and the *fcntl()* locks in IEEE Std 1003.1-2001 share a common goal:
 7426 increasing concurrency among readers, thus increasing throughput and decreasing delay.

7427 However, the read-write locks have two features not present in the *fcntl()* locks. First, under
 7428 priority scheduling, read-write locks are granted in priority order. Second, also under priority
 7429 scheduling, writer starvation is prevented by giving writers preference over readers of equal or
 7430 lower priority.

7431 Also, read-write locks can be used in systems lacking a file system, such as those conforming to
 7432 the minimal realtime system profile of IEEE Std 1003.13-1998.

7433 **History of Resolution Issues**

7434 Based upon some balloting objections, early drafts specified the behavior of threads waiting on a
 7435 read-write lock during the execution of a signal handler, as if the thread had not called the lock
 7436 operation. However, this specified behavior would require implementations to establish
 7437 internal signal handlers even though this situation would be rare, or never happen for many
 7438 programs. This would introduce an unacceptable performance hit in comparison to the little
 7439 additional functionality gained. Therefore, the behavior of read-write locks and signals was
 7440 reverted back to its previous mutex-like specification.

7441 **B.2.9.7** *Thread Interactions with Regular File Operations*

7442 There is no additional rationale provided for this section.

7443 **B.2.10** **Sockets**

7444 The base document for the sockets interfaces in IEEE Std 1003.1-2001 is the XNS, Issue 5.2
7445 specification. This was primarily chosen as it aligns with IPv6. Additional material has been
7446 added from IEEE Std 1003.1g-2000, notably socket concepts, raw sockets, the *pselect()* function,
7447 the *socketmark()* function, and the `<sys/select.h>` header.

7448 **B.2.10.1** *Address Families*

7449 There is no additional rationale provided for this section.

7450 **B.2.10.2** *Addressing*

7451 There is no additional rationale provided for this section.

7452 **B.2.10.3** *Protocols*

7453 There is no additional rationale provided for this section.

7454 **B.2.10.4** *Routing*

7455 There is no additional rationale provided for this section.

7456 **B.2.10.5** *Interfaces*

7457 There is no additional rationale provided for this section.

7458 **B.2.10.6** *Socket Types*

7459 The type `socklen_t` was invented to cover the range of implementations seen in the field. The
7460 intent of `socklen_t` is to be the type for all lengths that are naturally bounded in size; that is, that
7461 they are the length of a buffer which cannot sensibly become of massive size: network addresses,
7462 host names, string representations of these, ancillary data, control messages, and socket options
7463 are examples. Truly boundless sizes are represented by `size_t` as in *read()*, *write()*, and so on.

7464 All `socklen_t` types were originally (in BSD UNIX) of type `int`. During the development of
7465 IEEE Std 1003.1-2001, it was decided to change all buffer lengths to `size_t`, which appears at face
7466 value to make sense. When dual mode 32/64-bit systems came along, this choice unnecessarily
7467 complicated system interfaces because `size_t` (with `long`) was a different size under ILP32 and
7468 LP64 models. Reverting to `int` would have happened except that some implementations had
7469 already shipped 64-bit-only interfaces. The compromise was a type which could be defined to be
7470 any size by the implementation: `socklen_t`.

7471 **B.2.10.7** *Socket I/O Mode*

7472 There is no additional rationale provided for this section.

7473 *B.2.10.8 Socket Owner*

7474 There is no additional rationale provided for this section.

7475 *B.2.10.9 Socket Queue Limits*

7476 There is no additional rationale provided for this section.

7477 *B.2.10.10 Pending Error*

7478 There is no additional rationale provided for this section.

7479 *B.2.10.11 Socket Receive Queue*

7480 There is no additional rationale provided for this section.

7481 *B.2.10.12 Socket Out-of-Band Data State*

7482 There is no additional rationale provided for this section.

7483 *B.2.10.13 Connection Indication Queue*

7484 There is no additional rationale provided for this section.

7485 *B.2.10.14 Signals*

7486 There is no additional rationale provided for this section.

7487 *B.2.10.15 Asynchronous Errors*

7488 There is no additional rationale provided for this section.

7489 *B.2.10.16 Use of Options*

7490 There is no additional rationale provided for this section.

7491 *B.2.10.17 Use of Sockets for Local UNIX Connections*

7492 There is no additional rationale provided for this section.

7493 *B.2.10.18 Use of Sockets over Internet Protocols*

7494 A raw socket allows privileged users direct access to a protocol; for example, raw access to the
7495 IP and ICMP protocols is possible through raw sockets. Raw sockets are intended for
7496 knowledgeable applications that wish to take advantage of some protocol feature not directly
7497 accessible through the other sockets interfaces.

7498 *B.2.10.19 Use of Sockets over Internet Protocols Based on IPv4*

7499 There is no additional rationale provided for this section.

7500 *B.2.10.20 Use of Sockets over Internet Protocols Based on IPv6*

7501 The Open Group Base Resolution bwg2001-012 is applied, clarifying that IPv6 implementations
7502 are required to support use of AF_INET6 sockets over IPv4.

7503 **B.2.11 Tracing**

7504 The organization of the tracing rationale differs from the traditional rationale in that this tracing
 7505 rationale text is written against the trace interface as a whole, rather than against the individual
 7506 components of the trace interface or the normative section in which those components are
 7507 defined. Therefore the sections below do not parallel the sections of normative text in
 7508 IEEE Std 1003.1-2001.

7509 *B.2.11.1 Objectives*

7510 The intended uses of tracing are application-system debugging during system development, as a
 7511 “flight recorder” for maintenance of fielded systems, and as a performance measurement tool. In
 7512 all of these intended uses, the vendor-supplied computer system and its software are, for this
 7513 discussion, assumed error-free; the intent being to debug the user-written and/or third-party
 7514 application code, and their interactions. Clearly, problems with the vendor-supplied system and
 7515 its software will be uncovered from time to time, but this is a byproduct of the primary activity,
 7516 debugging user code.

7517 Another need for defining a trace interface in POSIX stems from the objective to provide an
 7518 efficient portable way to perform benchmarks. Existing practice shows that such interfaces are
 7519 commonly used in a variety of systems but with little commonality. As part of the benchmarking
 7520 needs, two aspects within the trace interface must be considered.

7521 The first, and perhaps more important one, is the qualitative aspect.

7522 The second is the quantitative aspect.

7523 • Qualitative Aspect

7524 To better understand this aspect, let us consider an example. Suppose that you want to
 7525 organize a number of actions to be performed during the day. Some of these actions are
 7526 known at the beginning of the day. Some others, which may be more or less important, will
 7527 be triggered by reading your mail. During the day you will make some phone calls and
 7528 synchronously receive some more information. Finally you will receive asynchronous phone
 7529 calls that also will trigger actions. If you, or somebody else, examines your day at work, you,
 7530 or he, can discover that you have not efficiently organized your work. For instance, relative
 7531 to the phone calls you made, would it be preferable to make some of these early in the
 7532 morning? Or to delay some others until the end of the day? Relative to the phone calls you
 7533 have received, you might find that somebody you called in the morning has called you 10
 7534 times while you were performing some important work. To examine, afterwards, your day at
 7535 work, you record in sequence all the trace events relative to your work. This should give you
 7536 a chance of organizing your next day at work.

7537 This is the qualitative aspect of the trace interface. The user of a system needs to keep a trace
 7538 of particular points the application passes through, so that he can eventually make some
 7539 changes in the application and/or system configuration, to give the application a chance of
 7540 running more efficiently.

7541 • Quantitative Aspect

7542 This aspect concerns primarily realtime applications, where missed deadlines can be
 7543 undesirable. Although there are, in IEEE Std 1003.1-2001, some interfaces useful for such
 7544 applications (timeouts, execution time monitoring, and so on), there are no APIs to aid in the
 7545 tuning of a realtime application’s behavior (**timespec** in timeouts, length of message queues,
 7546 duration of driver interrupt service routine, and so on). The tuning of an application needs a
 7547 means of recording timestamped important trace events during execution in order to analyze
 7548 offline, and eventually, to tune some realtime features (redesign the system with less

7549 functionalities, readjust timeouts, redesign driver interrupts, and so on).

7550 Detailed Objectives

7551 Objectives were defined to build the trace interface and are kept for historical interest. Although
7552 some objectives are not fully respected in this trace interface, the concept of the POSIX trace
7553 interface assumes the following points:

- 7554 1. It must be possible to trace both system and user trace events concurrently.
- 7555 2. It must be possible to trace per-process trace events and also to trace system trace events
7556 which are unrelated to any particular process. A per-process trace event is either user-
7557 initiated or system-initiated.
- 7558 3. It must be possible to control tracing on a per-process basis from either inside or outside
7559 the process.
- 7560 4. It must be possible to control tracing on a per-thread basis from inside the enclosing
7561 process.
- 7562 5. Trace points must be controllable by trace event type ID from inside and outside of the
7563 process. Multiple trace points can have the same trace event type ID, and will be controlled
7564 jointly.
- 7565 6. Recording of trace events is dependent on both trace event type ID and the
7566 process/thread. Both must be enabled in order to record trace events. System trace events
7567 may or may not be handled differently.
- 7568 7. The API must not mandate the ability to control tracing for more than one process at the
7569 same time.
- 7570 8. There is no objective for trace control on anything bigger than a process; for example,
7571 group or session.
- 7572 9. Trace propagation and control:
 - 7573 a. Trace propagation across *fork()* is optional; the default is to not trace a child process.
 - 7574 b. Trace control must span *pthread_create()* operations; that is, if a process is being
7575 traced, any thread will be traced as well if this thread allows tracing. The default is to
7576 allow tracing.
- 7577 10. Trace control must not span *exec* or *posix_spawn()* operations.
- 7578 11. A triggering API is not required. The triggering API is the ability to command or stop
7579 tracing based on the occurrence of a specific trace event other than a
7580 POSIX_TRACE_START trace event or a POSIX_TRACE_STOP trace event.
- 7581 12. Trace log entries must have timestamps of implementation-defined resolution.
7582 Implementations are exhorted to support at least microsecond resolution. When a trace log
7583 entry is retrieved, it must have timestamp, PC address, PID, and TID of the entity that
7584 generated the trace event.
- 7585 13. Independently developed code should be able to use trace facilities without coordination
7586 and without conflict.
- 7587 14. Even if the trace points in the trace calls are not unique, the trace log entries (after any
7588 processing) must be uniquely identified as to trace point.
- 7589 15. There must be a standard API to read the trace stream.

- 7590 16. The format of the trace stream and the trace log is opaque and unspecified.
- 7591 17. It must be possible to read a completed trace, if recorded on some suitable non-volatile
7592 storage, even subsequent to a power cycle or subsequent cold boot of the system.
- 7593 18. Support of analysis of a trace log while it is being formed is implementation-defined.
- 7594 19. The API must allow the application to write trace stream identification information into
7595 the trace stream and to be able to retrieve it, without it being overwritten by trace entries,
7596 even if the trace stream is full.
- 7597 20. It must be possible to specify the destination of trace data produced by trace events.
- 7598 21. It must be possible to have different trace streams, and for the tracing enabled by one trace
7599 stream to be completely independent of the tracing of another trace stream.
- 7600 22. It must be possible to trace events from threads in different CPUs.
- 7601 23. The API must support one or more trace streams per-system, and one or more trace
7602 streams per-process, up to an implementation-defined set of per-system and per-process
7603 maximums.
- 7604 24. It must be possible to determine the order in which the trace events happened, without
7605 necessarily depending on the clock, up to an implementation-defined time resolution.
- 7606 25. For performance reasons, the trace event point call(s) must be implementable as a macro
7607 (see the ISO POSIX-1: 1996 standard, 1.3.4, Statement 2).
- 7608 26. IEEE Std 1003.1-2001 must not define the trace points which a conforming system must
7609 implement, except for trace points used in the control of tracing.
- 7610 27. The APIs must be thread-safe, and trace points should be lock-free (that is, not require a
7611 lock to gain exclusive access to some resource).
- 7612 28. The user-provided information associated with a trace event is variable-sized, up to some
7613 maximum size.
- 7614 29. Bounds on record and trace stream sizes:
- 7615 a. The API must permit the application to declare the upper bounds on the length of an
7616 application data record. The system must return the limit it used. The limit used may
7617 be smaller than requested.
- 7618 b. The API must permit the application to declare the upper bounds on the size of trace
7619 streams. The system must return the limit it used. The limit used may be different,
7620 either larger or smaller, than requested.
- 7621 30. The API must be able to pass any fundamental data type, and a structured data type
7622 composed only of fundamental types. The API must be able to pass data by reference,
7623 given only as an address and a length. Fundamental types are the POSIX.1 types (see the
7624 `<sys/types.h>` header) plus those defined in the ISO C standard.
- 7625 31. The API must apply the POSIX notions of ownership and permission to recorded trace
7626 data, corresponding to the sources of that data.

7627 **Comments on Objectives**

7628 **Note:** In the following comments, numbers in square brackets refer to the above objectives.

7629 It is necessary to be able to obtain a trace stream for a complete activity. Thus there is a
7630 requirement to be able to trace both application and system trace events. A per-process trace
7631 event is either user-initiated, like the *write()* function, or system-initiated, like a timer expiration.
7632 There is also a need to be able to trace an entire process' activity even when it has threads in
7633 multiple CPUs. To avoid excess trace activity, it is necessary to be able to control tracing on a
7634 trace event type basis.

7635 [Objectives 1,2,5,22]

7636 There is a need to be able to control tracing on a per-process basis, both from inside and outside
7637 the process; that is, a process can start a trace activity on itself or any other process. There is also
7638 the perceived need to allow the definition of a maximum number of trace streams per system.

7639 [Objectives 3,23]

7640 From within a process, it is necessary to be able to control tracing on a per-thread basis. This
7641 provides an additional filtering capability to keep the amount of traced data to a minimum. It
7642 also allows for less ambiguity as to the origin of trace events. It is recognized that thread-level
7643 control is only valid from within the process itself. It is also desirable to know the maximum
7644 number of trace streams per process that can be started. The API should not require thread
7645 synchronization or mandate priority inversions that would cause the thread to block. However,
7646 the API must be thread-safe.

7647 [Objectives 4,23,24,27]

7648 There was no perceived objective to control tracing on anything larger than a process; for
7649 example, a group or session. Also, the ability to start or stop a trace activity on multiple
7650 processes atomically may be very difficult or cumbersome in some implementations.

7651 [Objectives 6,8]

7652 It is also necessary to be able to control tracing by trace event type identifier, sometimes called a
7653 trace hook ID. However, there is no mandated set of system trace events, since such trace points
7654 are implementation-defined. The API must not require from the operating system facilities that
7655 are not standard.

7656 [Objectives 6,26]

7657 Trace control must span *fork()* and *pthread_create()*. If not, there will be no way to ensure that an
7658 application's activity is entirely traced. The newly forked child would not be able to turn on its
7659 tracing until after it obtained control after the fork, and trace control externally would be even
7660 more problematic.

7661 [Objective 9]

7662 Since *exec* and *posix_spawn()* represent a complete change in the execution of a task (a new
7663 program), trace control need not persist over an *exec* or *posix_spawn()*.

7664 [Objective 10]

7665 Where trace activities are started on multiple processes, these trace activities should not interfere
7666 with each other.

7667 [Objective 21]

7668 There is no need for a triggering objective, primarily for performance reasons; see also Section
7669 B.2.11.8 (on page 203), rationale on triggering.

7670 [Objective 11]

7671 It must be possible to determine the origin of each traced event. The process and thread
7672 identifiers for each trace event are needed. Also there was a perceived need for a user-specifiable
7673 origin, but it was felt that this would create too much overhead.

- 7674 [Objectives 12,14]
- 7675 An allowance must be made for trace points to come embedded in software components from
7676 several different sources and vendors without requiring coordination.
7677 [Objective 13]
- 7678 There is a requirement to be able to uniquely identify trace points that may have the same trace
7679 stream identifier. This is only necessary when a trace report is produced.
7680 [Objectives 12,14]
- 7681 Tracing is a very performance-sensitive activity, and will therefore likely be implemented at a
7682 low level within the system. Hence the interface must not mandate any particular buffering or
7683 storage method. Therefore, a standard API is needed to read a trace stream. Also the interface
7684 must not mandate the format of the trace data, and the interface must not assume a trace storage
7685 method. Due to the possibility of a monolithic kernel and the possible presence of multiple
7686 processes capable of running trace activities, the two kinds of trace events may be stored in two
7687 separate streams for performance reasons. A mandatory dump mechanism, common in some
7688 existing practice, has been avoided to allow the implementation of this set of functions on small
7689 realtime profiles for which the concept of a file system is not defined. The trace API calls should
7690 be implemented as macros.
7691 [Objectives 15,16,25,30]
- 7692 Since a trace facility is a valuable service tool, the output (or log) of a completed trace stream
7693 that is written to permanent storage must be readable on other systems of the type that
7694 produced the trace log. Note that there is no objective to be able to interpret a trace log that was
7695 not successfully completed.
7696 [Objectives 17,18,19]
- 7697 For trace streams written to permanent storage, a way to specify the destination of the trace
7698 stream is needed.
7699 [Objective 20]
- 7700 There is a requirement to be able to depend on the ordering of trace events up to some
7701 implementation-defined time interval. For example, there is a need to know the time period
7702 during which, if trace events are closer together, their ordering is unspecified. Events that occur
7703 within an interval smaller than this resolution may or may not be read back in the correct order.
7704 [Objective 24]
- 7705 The application should be able to know how much data can be traced. When trace event types
7706 can be filtered, the application should be able to specify the approximate maximum amount of
7707 data that will be traced in a trace event so resources can be more efficiently allocated.
7708 [Objectives 28,29]
- 7709 Users should not be able to trace data to which they would not normally have access. System
7710 trace events corresponding to a process/thread should be associated with the ownership of that
7711 process/thread.
7712 [Objective 31]

7713 B.2.11.2 Trace Model

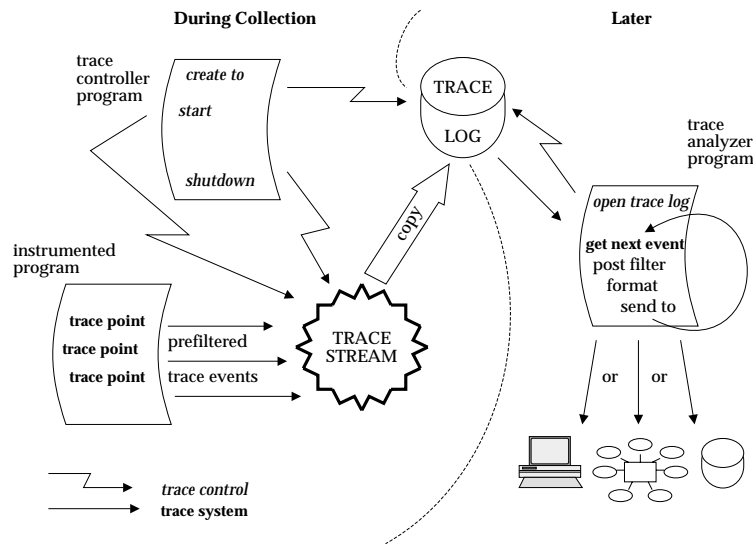
7714 **Introduction**

7715 The model is based on two base entities: the “Trace Stream” and the “Trace Log”, and a recorded
 7716 unit called the “Trace Event”. The possibility of using Trace Streams and Trace Logs separately
 7717 gives two use dimensions and solves both the performance issue and the full-information
 7718 system issue. In the case of a trace stream without log, specific information, although reduced in
 7719 quantity, is required to be registered, in a possibly small realtime system, with as little overhead
 7720 as possible. The Trace Log option has been added for small realtime systems. In the case of a
 7721 trace stream with log, considerable complex application-specific information needs to be
 7722 collected.

7723 **Trace Model Description**

7724 The trace model can be examined for three different subfunctions: Application Instrumentation,
 7725 Trace Operation Control, and Trace Analysis.

7726



7727 **Figure B-2** Trace System Overview: for Offline Analysis

7728 Each of these subfunctions requires specific characteristics of the trace mechanism API.

7729 • Application Instrumentation

7730 When instrumenting an application, the programmer is not concerned about the future use of
 7731 the trace events in the trace stream or the trace log, the full policy of the trace stream, or the
 7732 eventual pre-filtering of trace events. But he is concerned about the correct determination of
 7733 the specific trace event type identifier, regardless of how many independent libraries are
 7734 used in the same user application; see Figure B-2 and Figure B-3 (on page 186).

7735 This trace API provides the necessary operations to accomplish this subfunction. This is done
 7736 by providing functions to associate a programmer-defined name with an implementation-
 7737 defined trace event type identifier (see the *posix_trace_eventid_open()* function), and to send
 7738 this trace event into a potential trace stream (see the *posix_trace_event()* function).

7739 • Trace Operation Control

7740 When controlling the recording of trace events in a trace stream, the programmer is
 7741 concerned with the correct initialization of the trace mechanism (that is, the sizing of the
 7742 trace stream), the correct retention of trace events in a permanent storage, the correct
 7743 dynamic recording of trace events, and so on.

7744 This trace API provides the necessary material to permit this efficiently. This is done by
 7745 providing functions to initialize a new trace stream, and optionally a trace log:

- 7746 — Trace Stream Attributes Object Initialization (see *posix_trace_attr_init()*)
- 7747 — Functions to Retrieve or Set Information About a Trace Stream (see
 7748 *posix_trace_attr_getgenversion()*)
- 7749 — Functions to Retrieve or Set the Behavior of a Trace Stream (see
 7750 *posix_trace_attr_getinherited()*)
- 7751 — Functions to Retrieve or Set Trace Stream Size Attributes (see
 7752 *posix_trace_attr_getmaxuseventsize()*)
- 7753 — Trace Stream Initialization, Flush, and Shutdown from a Process (see *posix_trace_create()*)
- 7754 — Clear Trace Stream and Trace Log (see *posix_trace_clear()*)

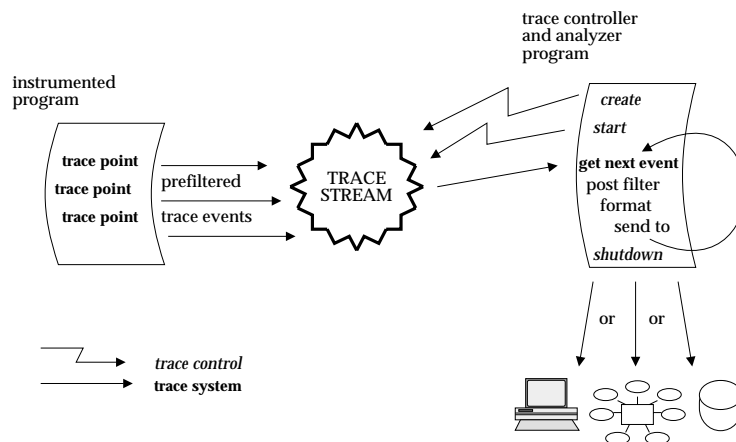
7755 To select the trace event types that are to be traced:

- 7756 — Manipulate Trace Event Type Identifier (see *posix_trace_trid_eventid_open()*)
- 7757 — Iterate over a Mapping of Trace Event Type (see *posix_trace_eventtypelist_getnext_id()*)
- 7758 — Manipulate Trace Event Type Sets (see *posix_trace_eventset_empty()*)
- 7759 — Set Filter of an Initialized Trace Stream (see *posix_trace_set_filter()*)

7760 To control the execution of an active trace stream:

- 7761 — Trace Start and Stop (see *posix_trace_start()*)
- 7762 — Functions to Retrieve the Trace Attributes or Trace Statuses (see *posix_trace_get_attr()*)

7763



7764 **Figure B-3** Trace System Overview: for Online Analysis

7765 • Trace Analysis

7766 Once correctly recorded, on permanent storage or not, an ultimate activity consists of the
 7767 analysis of the recorded information. If the recorded data is on permanent storage, a specific
 7768 open operation is required to associate a trace stream to a trace log.

7769 The first intent of the group was to request the presence of a system identification structure
 7770 in the trace stream attribute. This was, for the application, to allow some portable way to
 7771 process the recorded information. However, there is no requirement that the **utsname**
 7772 structure, on which this system identification was based, be portable from one machine to
 7773 another, so the contents of the attribute cannot be interpreted correctly by an application
 7774 conforming to IEEE Std 1003.1-2001.

7775 This modification has been incorporated and requests that some unspecified information be
 7776 recorded in the trace log in order to fail opening it if the analysis process and the controller
 7777 process were running in different types of machine, but does not request that this
 7778 information be accessible to the application. This modification has implied a modification in
 7779 the *posix_trace_open()* function error code returns.

7780 This trace API provides functions to:

- 7781 — Extract trace stream identification attributes (see *posix_trace_attr_getgenversion()*)
- 7782 — Extract trace stream behavior attributes (see *posix_trace_attr_getinherited()*)
- 7783 — Extract trace event, stream, and log size attributes (see
 7784 *posix_trace_attr_getmaxuseventsizesize()*)
- 7785 — Look up trace event type names (see *posix_trace_eventid_get_name()*)
- 7786 — Iterate over trace event type identifiers (see *posix_trace_eventtypelist_getnext_id()*)
- 7787 — Open, rewind, and close a trace log (see *posix_trace_open()*)
- 7788 — Read trace stream attributes and status (see *posix_trace_get_attr()*)
- 7789 — Read trace events (see *posix_trace_getnext_event()*)

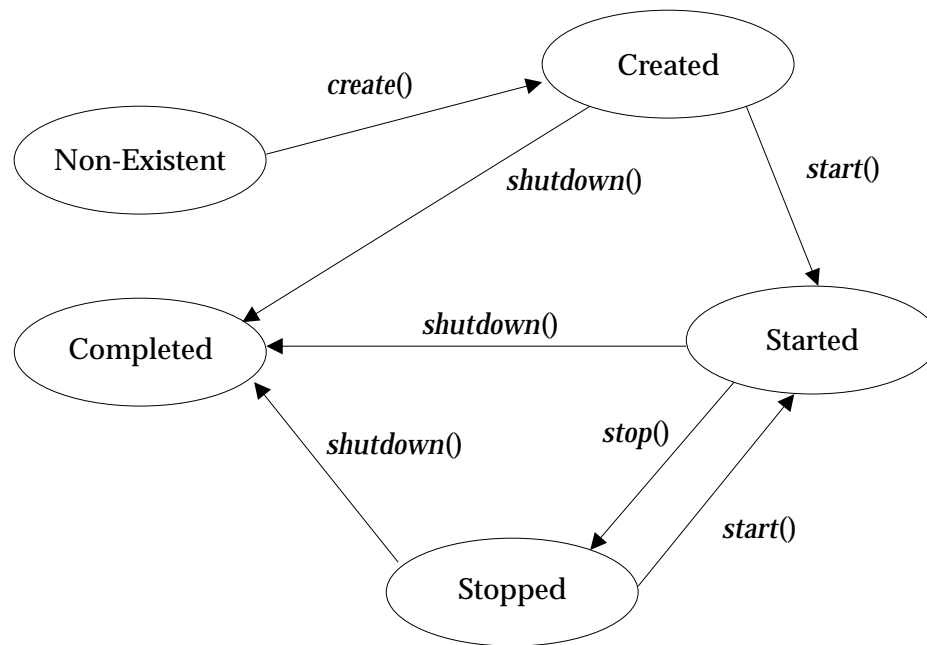
7790 Due to the following two reasons:

- 7791 1. The requirement that the trace system must not add unacceptable overhead to the traced
 7792 process and so that the trace event point execution must be fast
- 7793 2. The traced application does not care about tracing errors

7794 the trace system cannot return any internal error to the application. Internal error conditions can
 7795 range from unrecoverable errors that will force the active trace stream to abort, to small errors
 7796 that can affect the quality of tracing without aborting the trace stream. The group decided to
 7797 define a system trace event to report to the analysis process such internal errors. It is not the
 7798 intention of IEEE Std 1003.1-2001 to require an implementation to report an internal error that
 7799 corrupts or terminates tracing operation. The implementor is free to decide which internal
 7800 documented errors, if any, the trace system is able to report.

7801 **States of a Trace Stream**

7802



7803

Figure B-4 Trace System Overview: States of a Trace Stream

7804

7805

7806

7807

7808

7809

7810

7811

7812

Figure B-4 shows the different states an active trace stream passes through. After the *posix_trace_create()* function call, a trace stream becomes **CREATED** and a trace stream is associated for the future collection of trace events. The status of the trace stream is **POSIX_TRACE_SUSPENDED**. The state becomes **STARTED** after a call to the *posix_trace_start()* function, and the status becomes **POSIX_TRACE_RUNNING**. In this state, all trace events that are not filtered out will be stored into the trace stream. After a call to *posix_trace_stop()*, the trace stream becomes **STOPPED** (and the status **POSIX_TRACE_SUSPENDED**). In this state, no new trace events will be recorded in the trace stream, but previously recorded trace events may continue to be read.

7813

7814

7815

7816

7817

After a call to *posix_trace_shutdown()*, the trace stream is in the state **COMPLETED**. The trace stream no longer exists but, if the Trace Log option is supported, all the information contained in it has been logged. If a log object has not been associated with the trace stream at the creation, it is the responsibility of the trace controller process to not shut the trace stream down while trace events remain to be read in the stream.

7818

Tracing All Processes

7819

7820

7821

7822

Some implementations have a tracing subsystem with the ability to trace all processes. This is useful to debug some types of device drivers such as those for ATM or X25 adapters. These types of adapters are used by several independent processes, that are not issued from the same process.

7823

7824

7825

7826

The POSIX trace interface does not define any constant or option to create a trace stream tracing all processes. POSIX.1 does not prevent this type of implementation and an implementor is free to add this capability. Nevertheless, the trace interface allows tracing of all the system trace events and all the processes issued from the same process.

7827 If such a tracing system capability has to be implemented, when a trace stream is created, it is
7828 recommended that a constant named `POSIX_TRACE_ALLPROC` be used instead of the process
7829 identifier in the argument of the `posix_trace_create()` or `posix_trace_create_withlog()` function. A
7830 possible value for `POSIX_TRACE_ALLPROC` may be `-1` instead of a real process identifier.

7831 The implementor has to be aware that there is some impact on the tracing behavior as defined in
7832 the POSIX trace interface. For example:

- 7833 • If the default value for the inheritance attribute is set to
7834 `POSIX_TRACE_CLOSE_FOR_CHILD`, the implementation has to stop tracing for the child
7835 process.
- 7836 • The trace controller which is creating this type of trace stream must have the appropriate
7837 privilege to trace all the processes.

7838 Trace Storage

7839 The model is based on two types of trace events: system trace events and user-defined trace
7840 events. The internal representation of trace events is implementation-defined, and so the
7841 implementor is free to choose the more suitable, practical, and efficient way to design the
7842 internal management of trace events. For the timestamping operation, the model does not
7843 impose the `CLOCK_REALTIME` or any other clock. The buffering allocation and operation
7844 follow the same principle. The implementor is free to use one or more buffers to record trace
7845 events; the interface assumes only a logical trace stream of sequentially recorded trace events.
7846 Regarding flushing of trace events, the interface allows the definition of a trace log object which
7847 typically can be a file. But the group was also aware of defining functions to permit the use of
7848 this interface in small realtime systems, which may not have general file system capabilities. For
7849 instance, the three functions `posix_trace_getnext_event()` (blocking),
7850 `posix_trace_timedgetnext_event()` (blocking with timeout), and `posix_trace_trygetnext_event()`
7851 (non-blocking) are proposed to read the recorded trace events.

7852 The policy to be used when the trace stream becomes full also relies on common practice:

- 7853 • For an active trace stream, the `POSIX_TRACE_LOOP` trace stream policy permits automatic
7854 overrun (overwrite of oldest trace events) while waiting for some user-defined condition to
7855 cause tracing to stop. By contrast, the `POSIX_TRACE_UNTIL_FULL` trace stream policy
7856 requires the system to stop tracing when the trace stream is full. However, if the trace stream
7857 that is full is at least partially emptied by a call to the `posix_trace_flush()` function or by calls
7858 to the `posix_trace_getnext_event()` function, the trace system will automatically resume
7859 tracing.

7860 If the Trace Log option is supported, the operation of the `POSIX_TRACE_FLUSH` policy is an
7861 extension of the `POSIX_TRACE_UNTIL_FULL` policy. The automatic free operation (by
7862 flushing to the associated trace log) is added.

- 7863 • If a log is associated with the trace stream and this log is a regular file, these policies also
7864 apply for the log. One more policy, `POSIX_TRACE_APPEND`, is defined to allow indefinite
7865 extension of the log. Since the log destination can be any device or pseudo-device, the
7866 implementation may not be able to manipulate the destination as required by
7867 IEEE Std 1003.1-2001. For this reason, the behavior of the log full policy may be unspecified
7868 depending on the trace log type.

7869 The current trace interface does not define a service to preallocate space for a trace log file,
7870 because this space can be preallocated by means of a call to the `posix_fallocate()` function. This
7871 function could be called after the file has been opened, but before the trace stream is created.
7872 The `posix_fallocate()` function ensures that any required storage for regular file data is
7873 allocated on the file system storage media. If `posix_fallocate()` returns successfully,

7874 subsequent writes to the specified file data will not fail due to the lack of free space on the file
 7875 system storage media. Besides trace events, a trace stream also includes trace attributes and
 7876 the mapping from trace event names to trace event type identifiers. The implementor is free
 7877 to choose how to store the trace attributes and the trace event type map, but must ensure that
 7878 this information is not lost when a trace stream overrun occurs.

7879 B.2.11.3 Trace Programming Examples

7880 Several programming examples are presented to show the code of the different possible
 7881 subfunctions using a trace subsystem. All these programs need to include the <trace.h> header.
 7882 In the examples shown, error checking is omitted for more simplicity.

7883 Trace Operation Control

7884 These examples show the creation of a trace stream for another process; one which is already
 7885 trace instrumented. All the default trace stream attributes are used to simplify programming in
 7886 the first example. The second example shows more possibilities.

7887 First Example

```

7888 /* Caution. Error checks omitted */
7889 {
7890     trace_attr_t attr;
7891     pid_t pid = traced_process_pid;
7892     int fd;
7893     trace_id_t trid;
7894
7895     - - - - -
7896     /* Initialize trace stream attributes */
7897     posix_trace_attr_init(&attr);
7898     /* Open a trace log */
7899     fd=open("/tmp/mytracelog",...);
7900     /*
7901      * Create a new trace associated with a log
7902      * and with default attributes
7903      */
7904     posix_trace_create_withlog(pid, &attr, fd, &trid);
7905
7906     /* Trace attribute structure can now be destroyed */
7907     posix_trace_attr_destroy(&attr);
7908     /* Start of trace event recording */
7909     posix_trace_start(trid);
7910
7911     - - - - -
7912     /* Duration of tracing */
7913     - - - - -
7914     /* Stop and shutdown of trace activity */
7915     posix_trace_shutdown(trid);
7916     - - - - -
7917 }
```


7917 **Second Example**

7918 Between the initialization of the trace stream attributes and the creation of the trace stream,
 7919 these trace stream attributes may be modified; see **Trace Stream Attribute Manipulation** (on
 7920 page 195) for a specific programming example. Between the creation and the start of the trace
 7921 stream, the event filter may be set; after the trace stream is started, the event filter may be
 7922 changed. The setting of an event set and the change of a filter is shown in **Create a Trace Event**
 7923 **Type Set and Change the Trace Event Type Filter** (on page 195).

```

7924 /* Caution. Error checks omitted */
7925 {
7926     trace_attr_t attr;
7927     pid_t pid = traced_process_pid;
7928     int fd;
7929     trace_id_t trid;
7930     - - - - -
7931     /* Initialize trace stream attributes */
7932     posix_trace_attr_init(&attr);
7933     /* Attr default may be changed at this place; see example */
7934     - - - - -
7935     /* Create and open a trace log with R/W user access */
7936     fd=open("/tmp/mytracelog",O_WRONLY|O_CREAT,S_IRUSR|S_IWUSR);
7937     /* Create a new trace associated with a log */
7938     posix_trace_create_withlog(pid, &attr, fd, &trid);
7939     /*
7940      * If the Trace Filter option is supported
7941      * trace event type filter default may be changed at this place;
7942      * see example about changing the trace event type filter
7943      */
7944     posix_trace_start(trid);
7945     - - - - -
7946     /*
7947      * If you have an uninteresting part of the application
7948      * you can stop temporarily.
7949      *
7950      * posix_trace_stop(trid);
7951      * - - - - -
7952      * - - - - -
7953      * posix_trace_start(trid);
7954      */
7955     - - - - -
7956     /*
7957      * If the Trace Filter option is supported
7958      * the current trace event type filter can be changed
7959      * at any time (see example about how to set
7960      * a trace event type filter)
7961      */
7962     - - - - -
7963     /* Stop the recording of trace events */
7964     posix_trace_stop(trid);
7965     /* Shutdown the trace stream */
7966     posix_trace_shutdown(trid);
  
```

```

7967     /*
7968     * Destroy trace stream attributes; attr structure may have
7969     * been used during tracing to fetch the attributes
7970     */
7971     posix_trace_attr_destroy(&attr);
7972     - - - - -
7973 }

```

7974 Application Instrumentation

7975 This example shows an instrumented application. The code is included in a block of instructions,
 7976 perhaps a function from a library. Possibly in an initialization part of the instrumented
 7977 application, two user trace events names are mapped to two trace event type identifiers
 7978 (function *posix_trace_eventid_open()*). Then two trace points are programmed.

```

7979 /* Caution. Error checks omitted */
7980 {
7981     trace_event_id_t eventid1, eventid2;
7982     - - - - -
7983     /* Initialization of two trace event type ids */
7984     posix_trace_eventid_open("my_first_event",&eventid1);
7985     posix_trace_eventid_open("my_second_event",&eventid2);
7986     - - - - -
7987     - - - - -
7988     - - - - -
7989     /* Trace point */
7990     posix_trace_event(eventid1,NULL,0);
7991     - - - - -
7992     /* Trace point */
7993     posix_trace_event(eventid2,NULL,0);
7994     - - - - -
7995 }

```

7996 Trace Analyzer

7997 This example shows the manipulation of a trace log resulting from the dumping of a completed
 7998 trace stream. All the default attributes are used to simplify programming, and data associated
 7999 with a trace event is not shown in the first example. The second example shows more
 8000 possibilities.

8001 First Example

```

8002 /* Caution. Error checks omitted */
8003 {
8004     int fd;
8005     trace_id_t trid;
8006     posix_trace_event_info trace_event;
8007     char trace_event_name[TRACE_EVENT_NAME_MAX];
8008     int return_value;
8009     size_t returndatasize;
8010     int lost_event_number;
8011     - - - - -

```

```

8012     /* Open an existing trace log */
8013     fd=open("/tmp/tracelog", O_RDONLY);
8014     /* Open a trace stream on the open log */
8015     posix_trace_open(fd, &trid);
8016     /* Read a trace event */
8017     posix_trace_getnext_event(trid, &trace_event,
8018         NULL, 0, &returndatasize,&return_value);

8019     /* Read and print all trace event names out in a loop */
8020     while (return_value == NULL)
8021     {
8022         /*
8023          * Get the name of the trace event associated
8024          * with trid trace ID
8025          */
8026         posix_trace_eventid_get_name(trid, trace_event.event_id,
8027             trace_event_name);
8028         /* Print the trace event name out */
8029         printf("%s\n",trace_event_name);
8030         /* Read a trace event */
8031         posix_trace_getnext_event(trid, &trace_event,
8032             NULL, 0, &returndatasize,&return_value);
8033     }

8034     /* Close the trace stream */
8035     posix_trace_close(trid);
8036     /* Close the trace log */
8037     close(fd);
8038 }

```

8039 **Second Example**

8040 The complete example includes the two other examples in **Retrieve Information from a Trace**
8041 **Log** (on page 196) and in **Retrieve the List of Trace Event Types Used in a Trace Log** (on page
8042 197). For example, the *maxdatasize* variable is set in **Retrieve the List of Trace Event Types Used**
8043 **in a Trace Log** (on page 197).

```

8044     /* Caution. Error checks omitted */
8045     {
8046         int fd;
8047         trace_id_t trid;
8048         posix_trace_event_info trace_event;
8049         char trace_event_name[TRACE_EVENT_NAME_MAX];
8050         char * data;
8051         size_t maxdatasize=1024, returndatasize;
8052         int return_value;
8053         - - - - -

8054         /* Open an existing trace log */
8055         fd=open("/tmp/tracelog", O_RDONLY);
8056         /* Open a trace stream on the open log */
8057         posix_trace_open( fd, &trid);
8058         /*
8059          * Retrieve information about the trace stream which

```

```

8060         * was dumped in this trace log (see example)
8061         */
8062     - - - - -
8063     /* Allocate a buffer for trace event data */
8064     data=(char *)malloc(maxdatasize);
8065     /*
8066     * Retrieve the list of trace events used in this
8067     * trace log (see example)
8068     */
8069     - - - - -
8070     /* Read and print all trace event names and data out in a loop */
8071     while (1)
8072     {
8073     posix_trace_getnext_event(trid, &trace_event,
8074         data, maxdatasize, &returndatasize,&return_value);
8075         if (return_value != NULL) break;
8076         /*
8077         * Get the name of the trace event type associated
8078         * with trid trace ID
8079         */
8080         posix_trace_eventid_get_name(trid, trace_event.event_id,
8081             trace_event_name);
8082         {
8083         int i;
8084
8085         /* Print the trace event name out */
8086         printf("%s: ", trace_event_name);
8087         /* Print the trace event data out */
8088         for (i=0; i<returndatasize, i++) printf("%02.2X",
8089             (unsigned char)data[i]);
8090         printf("\n");
8091         }
8092
8093     /* Close the trace stream */
8094     posix_trace_close(trid);
8095     /* The buffer data is deallocated */
8096     free(data);
8097     /* Now the file can be closed */
8098     close(fd);
8099 }

```

8099 **Several Programming Manipulations**

8100 The following examples show some typical sets of operations needed in some contexts.

8101 **Trace Stream Attribute Manipulation**

8102 This example shows the manipulation of a trace stream attribute object in order to change the
8103 default value provided by a previous *posix_trace_attr_init()* call.

```
8104 /* Caution. Error checks omitted */
8105 {
8106     trace_attr_t attr;
8107     size_t logsize=100000;
8108     - - - - -
8109     /* Initialize trace stream attributes */
8110     posix_trace_attr_init(&attr);
8111     /* Set the trace name in the attributes structure */
8112     posix_trace_attr_setname(&attr, "my_trace");
8113     /* Set the trace full policy */
8114     posix_trace_attr_setstreamfullpolicy(&attr, POSIX_TRACE_LOOP);
8115     /* Set the trace log size */
8116     posix_trace_attr_setlogsize(&attr, logsize);
8117     - - - - -
8118 }
```

8119 **Create a Trace Event Type Set and Change the Trace Event Type Filter**

8120 This example is valid only if the Trace Event Filter option is supported. This example shows the
8121 manipulation of a trace event type set in order to change the trace event type filter for an existing
8122 active trace stream, which may be just-created, running, or suspended. Some sets of trace event
8123 types are well-known, such as the set of trace event types not associated with a process, some
8124 trace event types are just-built trace event types for this trace stream; one trace event type is the
8125 predefined trace event error type which is deleted from the trace event type set.

```
8126 /* Caution. Error checks omitted */
8127 {
8128     trace_id_t trid = existing_trace;
8129     trace_event_set_t set;
8130     trace_event_id_t trace_event1, trace_event2;
8131     - - - - -
8132     /* Initialize to an empty set of trace event types */
8133     /* (not strictly required because posix_trace_event_set_fill() */
8134     /* will ignore the prior contents of the event set.) */
8135     posix_trace_eventset_emptyset(&set);
8136     /*
8137     * Fill the set with all system trace events
8138     * not associated with a process
8139     */
8140     posix_trace_eventset_fill(&set, POSIX_TRACE_WOPID_EVENTS);
8141     /*
8142     * Get the trace event type identifier of the known trace event name
8143     * my_first_event for the trid trace stream
8144     */
8145     posix_trace_trid_eventid_open(trid, "my_first_event", &trace_event1);
8146     /* Add the set with this trace event type identifier */
8147     posix_trace_eventset_add_event(trace_event1, &set);
8148     /*
```

```

8149     * Get the trace event type identifier of the known trace event name
8150     * my_second_event for the trid trace stream
8151     */
8152     posix_trace_trid_eventid_open(trid, "my_second_event", &trace_event2);
8153     /* Add the set with this trace event type identifier */
8154     posix_trace_eventset_add_event(trace_event2, &set);
8155     - - - - -
8156     /* Delete the system trace event POSIX_TRACE_ERROR from the set */
8157     posix_trace_eventset_del_event(POSIX_TRACE_ERROR, &set);
8158     - - - - -
8159
8160     /* Modify the trace stream filter making it equal to the new set */
8161     posix_trace_set_filter(trid, &set, POSIX_TRACE_SET_EVENTSET);
8162     - - - - -
8163     /*
8164     * Now trace_event1, trace_event2, and all system trace event types
8165     * not associated with a process, except for the POSIX_TRACE_ERROR
8166     * system trace event type, are filtered out of (not recorded in) the
8167     * existing trace stream.
8168     */
8169     }

```

8169 Retrieve Information from a Trace Log

8170 This example shows how to extract information from a trace log, the dump of a trace stream.
8171 This code:

```

8172     • Asks if the trace stream has lost trace events
8173     • Extracts the information about the version of the trace subsystem which generated this trace
8174     log
8175     • Retrieves the maximum size of trace event data; this may be used to dynamically allocate an
8176     array for extracting trace event data from the trace log without overflow
8177     /* Caution. Error checks omitted */
8178     {
8179         struct posix_trace_status_info statusinfo;
8180         trace_attr_t attr;
8181         trace_id_t trid = existing_trace;
8182         size_t maxdatasize;
8183         char genversion[TRACE_NAME_MAX];
8184         - - - - -
8185         /* Get the trace stream status */
8186         posix_trace_get_status(trid, &statusinfo);
8187         /* Detect an overrun condition */
8188         if (statusinfo.posix_stream_overrun_status == POSIX_TRACE_OVERRUN)
8189             printf("trace events have been lost\n");
8190
8191         /* Get attributes from the trid trace stream */
8192         posix_trace_get_attr(trid, &attr);
8193         /* Get the trace generation version from the attributes */
8194         posix_trace_attr_getgenversion(&attr, genversion);
8195         /* Print the trace generation version out */
8196         printf("Information about Trace Generator:%s\n",genversion);

```

```

8196     /* Get the trace event max data size from the attributes */
8197     posix_trace_attr_getmaxdatasize(&attr, &maxdatasize);
8198     /* Print the trace event max data size out */
8199     printf("Maximum size of associated data:%d\n",maxdatasize);
8200     /* Destroy the trace stream attributes */
8201     posix_trace_attr_destroy(&attr);
8202 }

```

8203 **Retrieve the List of Trace Event Types Used in a Trace Log**

8204 This example shows the retrieval of a trace stream's trace event type list. This operation may be
8205 very useful if you are interested only in tracking the type of trace events in a trace log.

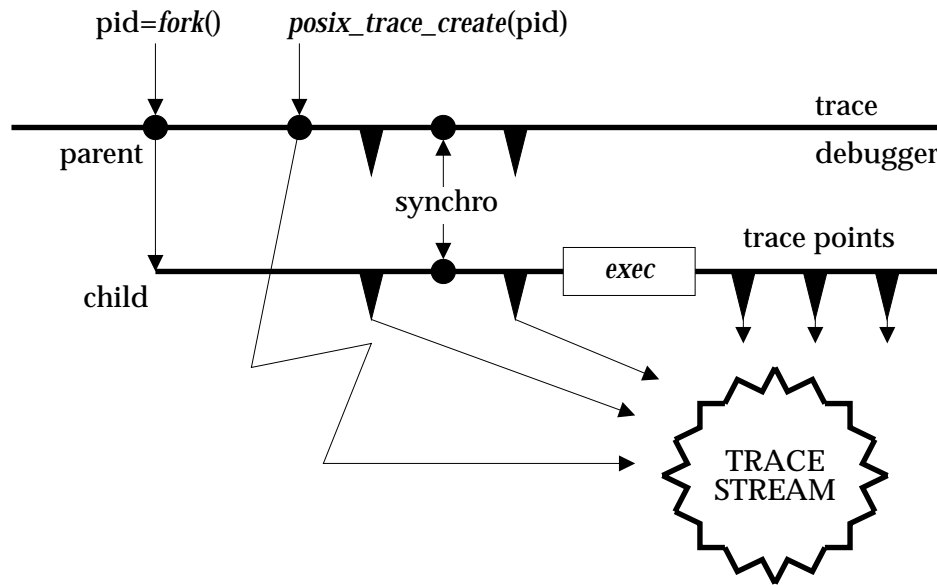
```

8206 /* Caution. Error checks omitted */
8207 {
8208     trace_id_t trid = existing_trace;
8209     trace_event_id_t event_id;
8210     char event_name[TRACE_EVENT_NAME_MAX];
8211     int return_value;
8212     - - - - -
8213
8214     /*
8215     * In a loop print all existing trace event names out
8216     * for the trid trace stream
8217     */
8217     while (1)
8218     {
8219         posix_trace_eventtypelist_getnext_id(trid, &event_id
8220             &return_value);
8221         if (return_value != NULL) break;
8222         /*
8223         * Get the name of the trace event associated
8224         * with trid trace ID
8225         */
8226         posix_trace_eventid_get_name(trid, event_id, event_name);
8227         /* Print the name out */
8228         printf("%s\n", event_name);
8229     }
8230 }

```

8231 B.2.11.4 Rationale on Trace for Debugging

8232



8233

Figure B-5 Trace Another Process

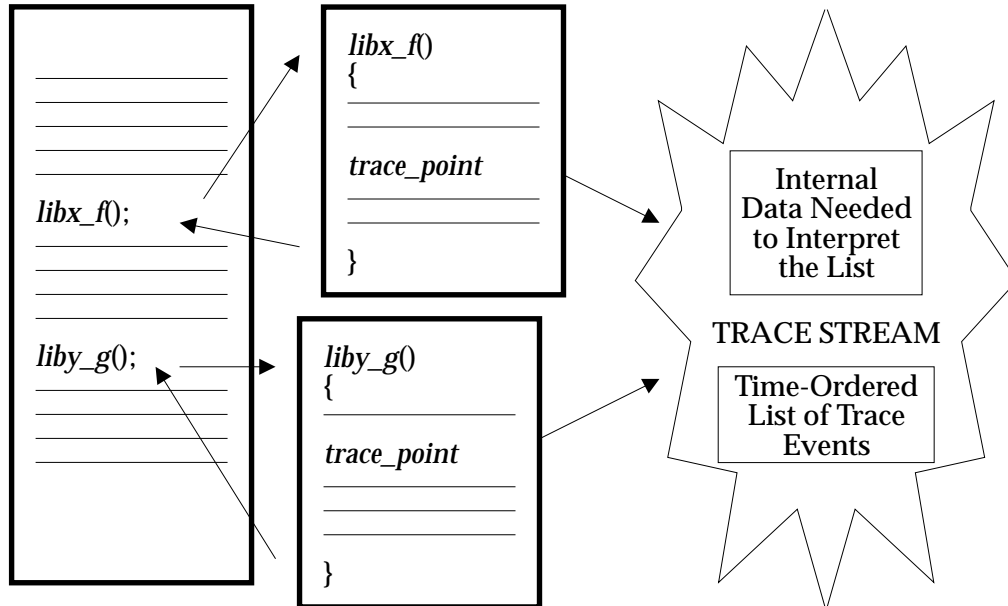
8234 Among the different possibilities offered by the trace interface defined in IEEE Std 1003.1-2001,
 8235 the debugging of an application is the most interesting one. Typical operations in the controlling
 8236 debugger process are to filter trace event types, to get trace events from the trace stream, to stop
 8237 the trace stream when the debugged process is executing uninteresting code, to start the trace
 8238 stream when some interesting point is reached, and so on. The interface defined in
 8239 IEEE Std 1003.1-2001 should define all the necessary base functions to allow this dynamic debug
 8240 handling.

8241 Figure B-5 shows an example in which the trace stream is created after the call to the `fork()`
 8242 function. If the user does not want to lose trace events, some synchronization mechanism
 8243 (represented in the figure) may be needed before calling the `exec` function, to give the parent a
 8244 chance to create the trace stream before the child begins the execution of its trace points.

8245 B.2.11.5 Rationale on Trace Event Type Name Space

8246 At first, the working group was in favor of the representation of a trace event type by an integer
 8247 (`event_name`). It seems that existing practice shows the weakness of such a representation. The
 8248 collision of trace event types is the main problem that cannot be simply resolved using this sort
 8249 of representation. Suppose, for example, that a third party designs an instrumented library. The
 8250 user does not have the source of this library and wants to trace his application which uses in
 8251 some part the third-party library. There is no means for him to know what are the trace event
 8252 types used in the instrumented library so he has some chance of duplicating some of them and
 8253 thus to obtain a contaminated tracing of his application.

8254



8255

Figure B-6 Trace Name Space Overview: With Third-Party Library

8256 There are requirements to allow program images containing pieces from various vendors to be
 8257 traced without also requiring those of any other vendors to coordinate their uses of the trace
 8258 facility, and especially the naming of their various trace event types and trace point IDs. The
 8259 chosen solution is to provide a very large name space, large enough so that the individual
 8260 vendors can give their trace types and tracepoint IDs sufficiently long and descriptive names
 8261 making the occurrence of collisions quite unlikely. The probability of collision is thus made
 8262 sufficiently low so that the problem may, as a practical matter, be ignored. By requirement, the
 8263 consequence of collisions will be a slight ambiguity in the trace streams; tracing will continue in
 8264 spite of collisions and ambiguities. “The show must go on”. The *posix_prog_address* member of
 8265 the **posix_trace_event_info** structure is used to allow trace streams to be unambiguously
 8266 interpreted, despite the fact that trace event types and trace event names need not be unique.

8267 The *posix_trace_eventid_open()* function is required to allow the instrumented third-party library
 8268 to get a valid trace event type identifier for its trace event names. This operation is, somehow,
 8269 an allocation, and the group was aware of proposing some deallocation mechanism which the
 8270 instrumented application could use to recover the resources used by a trace event type identifier.
 8271 This would have given the instrumented application the benefit of being capable of reusing a
 8272 possible minimum set of trace event type identifiers, but also the inconvenience to have,
 8273 possibly in the same trace stream, one trace event type identifier identifying two different trace
 8274 event types. After some discussions the group decided to not define such a function which
 8275 would make this API thicker for little benefit, the user having always the possibility of adding
 8276 identification information in the *data* member of the trace event structure.

8277 The set of the trace event type identifiers the controlling process wants to filter out is initialized
 8278 in the trace mechanism using the function *posix_trace_set_filter()*, setting the arguments
 8279 according to the definitions explained in *posix_trace_set_filter()*. This operation can be done
 8280 statically (when the trace is in the STOPPED state) or dynamically (when the trace is in the
 8281 STARTED state). The preparation of the filter is normally done using the function defined in
 8282 *posix_trace_eventtypelist_getnext_id()* and eventually the function
 8283 *posix_trace_eventtypelist_rewind()* in order to know (before the recording) the list of the potential

8284 set of trace event types that can be recorded. In the case of an active trace stream, this list may
8285 not be exhaustive. Actually, the target process may not have yet called the function
8286 *posix_trace_eventid_open()*. But it is a common practice, for a controlling process, to prepare the
8287 filtering of a future trace stream before its start. Therefore the user must have a way to get the
8288 trace event type identifier corresponding to a well-known trace event name before its future
8289 association by the pre-cited function. This is done by calling the *posix_trace_trid_eventid_open()*
8290 function, given the trace stream identifier and the trace name, and described hereafter. Because
8291 this trace event type identifier is associated with a trace stream identifier, where a unique
8292 process has initialized two or more traces, the implementation is expected to return the same
8293 trace event type identifier for successive calls to *posix_trace_trid_eventid_open()* with different
8294 trace stream identifiers. The *posix_trace_eventid_get_name()* function is used by the controller
8295 process to identify, by the name, the trace event type returned by a call to the
8296 *posix_trace_eventtypelist_getnext_id()* function.

8297 Afterwards, the set of trace event types is constructed using the functions defined in
8298 *posix_trace_eventset_empty()*, *posix_trace_eventset_fill()*, *posix_trace_eventset_add()*, and
8299 *posix_trace_eventset_del()*.

8300 A set of functions is provided devoted to the manipulation of the trace event type identifier and
8301 names for an active trace stream. All these functions require the trace stream identifier argument
8302 as the first parameter. The opacity of the trace event type identifier implies that the user cannot
8303 associate directly its well-known trace event name with the system-associated trace event type
8304 identifier.

8305 The *posix_trace_trid_eventid_open()* function allows the application to get the system trace event
8306 type identifier back from the system, given its well-known trace event name. This function is
8307 useful only when a controlling process needs to specify specific events to be filtered.

8308 The *posix_trace_eventid_get_name()* function allows the application to obtain a trace event name
8309 given its trace event type identifier. One possible use of this function is to identify the type of a
8310 trace event retrieved from the trace stream, and print it. The easiest way to implement this
8311 requirement, is to use a single trace event type map for all the processes whose maps are
8312 required to be identical. A more difficult way is to attempt to keep multiple maps identical at
8313 every call to *posix_trace_eventid_open()* and *posix_trace_trid_eventid_open()*.

8314 *B.2.11.6 Rationale on Trace Events Type Filtering*

8315 The most basic rationale for runtime and pre-registration filtering (selection/rejection) of trace
8316 event types is to prevent choking of the trace collection facility, and/or overloading of the
8317 computer system. Any worthwhile trace facility can bring even the largest computer to its
8318 knees. Otherwise, everything would be recorded and filtered after the fact; it would be much
8319 simpler, but impractical.

8320 To achieve debugging, measurement, or whatever the purpose of tracing, the filtering of trace
8321 event types is an important part of trace analysis. Due to the fact that the trace events are put
8322 into a trace stream and probably logged afterwards into a file, different levels of filtering—that
8323 is, rejection of trace event types—are possible.

8324 **Filtering of Trace Event Types Before Tracing**

8325 This function, represented by the `posix_trace_set_filter()` function in IEEE Std 1003.1-2001 (see
8326 `posix_trace_set_filter()`), selects, before or during tracing, the set of trace event types to be filtered
8327 out. It should be possible also (as OSF suggested in their ETAP trace specifications) to select the
8328 kernel trace event types to be traced in a system-wide fashion. These two functionalities are
8329 called the pre-filtering of trace event types.

8330 The restriction on the actual type used for the `trace_event_set_t` type is intended to guarantee
8331 that these objects can always be assigned, have their address taken, and be passed by value as
8332 parameters. It is not intended that this type be a structure including pointers to other data
8333 structures, as that could impact the portability of applications performing such operations. A
8334 reasonable implementation could be a structure containing an array of integer types.

8335 **Filtering of Trace Event Types at Runtime**

8336 It is possible to build this functionality using the `posix_trace_set_filter()` function. A privileged
8337 process or a privileged thread can get trace events from the trace stream of another process or
8338 thread, and thus specify the type of trace events to record into a file, using implementation-
8339 defined methods and interfaces. This functionality, called inline filtering of trace event types, is
8340 used for runtime analysis of trace streams.

8341 **Post-Mortem Filtering of Trace Event Types**

8342 The word “post-mortem” is used here to indicate that some unanticipated situation occurs
8343 during execution that does not permit a pre or inline filtering of trace events and that it is
8344 necessary to record all trace event types to have a chance to discover the problem afterwards.
8345 When the program stops, all the trace events recorded previously can be analyzed in order to
8346 find the solution. This functionality could be named the post-filtering of trace event types.

8347 **Discussions about Trace Event Type-Filtering**

8348 After long discussions with the parties involved in the process of defining the trace interface, it
8349 seems that the sensitivity to the filtering problem is different, but everybody agrees that the level
8350 of the overhead introduced during the tracing operation depends on the filtering method
8351 elected. If the time that it takes the trace event to be recorded can be neglected, the overhead
8352 introduced by the filtering process can be classified as follows:

8353 Pre-filtering System and process/thread-level overhead

8354 Inline-filtering Process/thread-level overhead

8355 Post-filtering No overhead; done offline

8356 The pre-filtering could be named “critical realtime” filtering in the sense that the filtering of
8357 trace event type is manageable at the user level so the user can lower to a minimum the filtering
8358 overhead at some user selected level of priority for the inline filtering, or delay the filtering to
8359 after execution for the post-filtering. The counterpart of this solution is that the size of the trace
8360 stream must be sufficient to record all the trace events. The advantage of the pre-filtering is that
8361 the utilization of the trace stream is optimized.

8362 Only pre-filtering is defined by IEEE Std 1003.1-2001. However, great care must be taken in
8363 specifying pre-filtering, so that it does not impose unacceptable overhead. Moreover, it is
8364 necessary to isolate all the functionality relative to the pre-filtering.

8365 The result of this rationale is to define a new option, the Trace Event Filter option, not
8366 necessarily implemented in small realtime systems, where system overhead is minimized to the
8367 extent possible.

8368 B.2.11.7 Tracing, pthread API

8369 The objective to be able to control tracing for individual threads may be in conflict with the
 8370 efficiency expected in threads with a *contentionscope* attribute of PTHREAD_SCOPE_PROCESS.
 8371 For these threads, context switches from one thread that has tracing enabled to another thread
 8372 that has tracing disabled may require a kernel call to inform the kernel whether it has to trace
 8373 system events executed by that thread or not. For this reason, it was proposed that the ability to
 8374 enable or disable tracing for PTHREAD_SCOPE_PROCESS threads be made optional, through
 8375 the introduction of a Trace Scope Process option. A trace implementation which did not
 8376 implement the Trace Scope Process option would not honor the tracing-state attribute of a
 8377 thread with PTHREAD_SCOPE_PROCESS; it would, however, honor the tracing-state attribute
 8378 of a thread with PTHREAD_SCOPE_SYSTEM. This proposal was rejected as:

- 8379 1. Removing desired functionality (per-thread trace control)
- 8380 2. Introducing counter-intuitive behavior for the tracing-state attribute
- 8381 3. Mixing logically orthogonal ideas (thread scheduling and thread tracing)
- 8382 [Objective 4]

8383 Finally, to solve this complex issue, this API does not provide *pthread_gettracingstate()*,
 8384 *pthread_settracingstate()*, *pthread_attr_gettracingstate()*, and *pthread_attr_settracingstate()*
 8385 interfaces. These interfaces force the thread implementation to add to the weight of the thread
 8386 and cause a revision of the threads libraries, just to support tracing. Worse yet,
 8387 *posix_trace_event()* must always test this per-thread variable even in the common case where it is
 8388 not used at all. Per-thread tracing is easy to implement using existing interfaces where necessary;
 8389 see the following example.

8390 **Example**

```
8391 /* Caution. Error checks omitted */
8392 static pthread_key_t my_key;
8393 static trace_event_id_t my_event_id;
8394 static pthread_once_t my_once = PTHREAD_ONCE_INIT;
8395
8396 void my_init(void)
8397 {
8398     (void) pthread_key_create(&my_key, NULL);
8399     (void) posix_trace_eventid_open("my", &my_event_id);
8400 }
8401
8402 int get_trace_flag(void)
8403 {
8404     pthread_once(&my_once, my_init);
8405     return (pthread_getspecific(my_key) != NULL);
8406 }
8407
8408 void set_trace_flag(int f)
8409 {
8410     pthread_once(&my_once, my_init);
8411     pthread_setspecific(my_key, f? &my_event_id: NULL);
8412 }
8413
8414 fn()
8415 {
8416     if (get_trace_flag())
8417         posix_trace_event(my_event_id, ...)
```

8414 }

8415 The above example does not implement third-party state setting.

8416 Lastly, per-thread tracing works poorly for threads with PTHREAD_SCOPE_PROCESS
8417 contention scope. These “library” threads have minimal interaction with the kernel and would
8418 have to explicitly set the attributes whenever they are context switched to a new kernel thread in
8419 order to trace system events. Such state was explicitly avoided in POSIX threads to keep
8420 PTHREAD_SCOPE_PROCESS threads lightweight.

8421 The reason that keeping PTHREAD_SCOPE_PROCESS threads lightweight is important is that
8422 such threads can be used not just for simple multi-processors but also for co-routine style
8423 programming (such as discrete event simulation) without inventing a new threads paradigm.
8424 Adding extra runtime cost to thread context switches will make using POSIX threads less
8425 attractive in these situations.

8426 *B.2.11.8 Rationale on Triggering*

8427 The ability to start or stop tracing based on the occurrence of specific trace event types has been
8428 proposed as a parallel to similar functionality appearing in logic analyzers. Such triggering, in
8429 order to be very useful, should be based not only on the trace event type, but on trace event-
8430 specific data, including tests of user-specified fields for matching or threshold values.

8431 Such a facility is unnecessary where the buffering of the stream is not a constraint, since such
8432 checks can be performed offline during post-mortem analysis.

8433 For example, a large system could incorporate a daemon utility to collect the trace records from
8434 memory buffers and spool them to secondary storage for later analysis. In the instances where
8435 resources are truly limited, such as embedded applications, the application incorporation of
8436 application code to test the circumstances of a trace event and call the trace point only if needed
8437 is usually straightforward.

8438 For performance reasons, the *posix_trace_event()* function should be implemented using a macro,
8439 so if the trace is inactive, the trace event point calls are latent code and must cost no more than a
8440 scalar test.

8441 The API proposed in IEEE Std 1003.1-2001 does not include any triggering functionality.

8442 *B.2.11.9 Rationale on Timestamp Clock*

8443 It has been suggested that the tracing mechanism should include the possibility of specifying the
8444 clock to be used in timestamping the trace events. When application trace events must be
8445 correlated to remote trace events, such a facility could provide a global time reference not
8446 available from a local clock. Further, the application may be driven by timers based on a clock
8447 different from that used for the timestamp, and the correlation of the trace to those untraced
8448 timer activities could be an important part of the analysis of the application.

8449 However, the tracing mechanism needs to be fast and just the provision of such an option can
8450 materially affect its performance. Leaving aside the performance costs of reading some clocks,
8451 this notion is also ill-defined when kernel trace events are to be traced by two applications
8452 making use of different tracing clocks. This can even happen within a single application where
8453 different parts of the application are served by different clocks. Another complication can occur
8454 when a clock is maintained strictly at the user level and is unavailable at the kernel level.

8455 It is felt that the benefits of a selectable trace clock do not match its costs. Applications that wish
8456 to correlate clocks other than the default tracing clock can include trace events with sample
8457 values of those other clocks, allowing correlation of timestamps from the various independent
8458 clocks. In any case, such a technique would be required when applications are sensitive to

8459 multiple clocks.

8460 *B.2.11.10 Rationale on Different Overrun Conditions*

8461 The analysis of the dynamic behavior of the trace mechanism shows that different overrun
8462 conditions may occur. The API must provide a means to manage such conditions in a portable
8463 way.

8464 **Overrun in Trace Streams Initialized with POSIX_TRACE_LOOP Policy**

8465 In this case, the user of the trace mechanism is interested in using the trace stream with
8466 POSIX_TRACE_LOOP policy to record trace events continuously, but ideally without losing any
8467 trace events. The online analyzer process must get the trace events at a mean speed equivalent to
8468 the recording speed. Should the trace stream become full, a trace stream overrun occurs. This
8469 condition is detected by getting the status of the active trace stream (function
8470 *posix_trace_get_status()*) and looking at the member *posix_stream_overrun_status* of the read
8471 **posix_stream_status** structure. In addition, two predefined trace event types are defined:

- 8472 1. The beginning of a trace overflow, to locate the beginning of an overflow when reading a
8473 trace stream
- 8474 2. The end of a trace overflow, to locate the end of an overflow, when reading a trace stream

8475 As a timestamp is associated with these predefined trace events, it is possible to know the
8476 duration of the overflow.

8477 **Overrun in Dumping Trace Streams into Trace Logs**

8478 The user lets the trace mechanism dump the trace stream initialized with
8479 POSIX_TRACE_FLUSH policy automatically into a trace log. If the dump operation is slower
8480 than the recording of trace events, the trace stream can overrun. This condition is detected by
8481 getting the status of the active trace stream (function *posix_trace_get_status()*) and looking at the
8482 member *posix_log_overrun_status* of the read **posix_stream_status** structure. This overrun
8483 indicates that the trace mechanism is not able to operate in this mode at this speed. It is the
8484 responsibility of the user to modify one of the trace parameters (the stream size or the trace
8485 event type filter, for instance) to avoid such overrun conditions, if overruns are to be prevented.
8486 The same already predefined trace event types (see **Overrun in Trace Streams Initialized with**
8487 **POSIX_TRACE_LOOP Policy**) are used to detect and to know the duration of an overflow.

8488 **Reading an Active Trace Stream**

8489 Although this trace API allows one to read an active trace stream with log while it is tracing, this
8490 feature can lead to false overflow origin interpretation: the trace log or the reader of the trace
8491 stream. Reading from an active trace stream with log is thus non-portable, and has been left
8492 unspecified.

8493 **B.2.12 Data Types**

8494 The requirement that additional types defined in this section end in “_t” was prompted by the
8495 problem of name space pollution. It is difficult to define a type (where that type is not one
8496 defined by IEEE Std 1003.1-2001) in one header file and use it in another without adding symbols
8497 to the name space of the program. To allow implementors to provide their own types, all
8498 conforming applications are required to avoid symbols ending in “_t”, which permits the
8499 implementor to provide additional types. Because a major use of types is in the definition of
8500 structure members, which can (and in many cases must) be added to the structures defined in
8501 IEEE Std 1003.1-2001, the need for additional types is compelling.

- 8502 The types, such as **ushort** and **ulong**, which are in common usage, are not defined in
 8503 IEEE Std 1003.1-2001 (although **ushort_t** would be permitted as an extension). They can be
 8504 added to `<sys/types.h>` using a feature test macro (see Section B.2.2.1 (on page 87)). A suggested
 8505 symbol for these is `_SYSIII`. Similarly, the types like **u_short** would probably be best controlled
 8506 by `_BSD`.
- 8507 Some of these symbols may appear in other headers; see Section B.2.2.2 (on page 88).
- 8508 **dev_t** This type may be made large enough to accommodate host-locality considerations
 8509 of networked systems.
- 8510 This type must be arithmetic. Earlier proposals allowed this to be non-arithmetic
 8511 (such as a structure) and provided a `samefile()` function for comparison.
- 8512 **gid_t** Some implementations had separated **gid_t** from **uid_t** before POSIX.1 was
 8513 completed. It would be difficult for them to coalesce them when it was
 8514 unnecessary. Additionally, it is quite possible that user IDs might be different from
 8515 group IDs because the user ID might wish to span a heterogeneous network,
 8516 where the group ID might not.
- 8517 For current implementations, the cost of having a separate **gid_t** will be only
 8518 lexical.
- 8519 **mode_t** This type was chosen so that implementations could choose the appropriate
 8520 integer type, and for compatibility with the ISO C standard. 4.3 BSD uses
 8521 **unsigned short** and the SVID uses **ushort**, which is the same. Historically, only the
 8522 low-order sixteen bits are significant.
- 8523 **nlink_t** This type was introduced in place of **short** for `st_nlink` (see the `<sys/stat.h>` header)
 8524 in response to an objection that **short** was too small.
- 8525 **off_t** This type is used only in `lseek()`, `fcntl()`, and `<sys/stat.h>`. Many implementations
 8526 would have difficulties if it were defined as anything other than **long**. Requiring
 8527 an integer type limits the capabilities of `lseek()` to four gigabytes. The ISO C
 8528 standard supplies routines that use larger types; see `fgetpos()` and `fsetpos()`. XSI-
 8529 conformant systems provide the `lseeko()` and `ftello()` functions that use larger
 8530 types.
- 8531 **pid_t** The inclusion of this symbol was controversial because it is tied to the issue of the
 8532 representation of a process ID as a number. From the point of view of a
 8533 conforming application, process IDs should be “magic cookies”¹ that are produced
 8534 by calls such as `fork()`, used by calls such as `waitpid()` or `kill()`, and not otherwise
 8535 analyzed (except that the sign is used as a flag for certain operations).
- 8536 The concept of a {PID_MAX} value interacted with this in early proposals. Treating
 8537 process IDs as an opaque type both removes the requirement for {PID_MAX} and
 8538 allows systems to be more flexible in providing process IDs that span a large range
 8539 of values, or a small one.
- 8540 Since the values in **uid_t**, **gid_t**, and **pid_t** will be numbers generally, and
 8541 potentially both large in magnitude and sparse, applications that are based on
- 8542 _____
- 8543 1. An historical term meaning: “An opaque object, or token, of determinate size, whose significance is known only to the entity
 8544 which created it. An entity receiving such a token from the generating entity may only make such use of the ‘cookie’ as is defined
 8545 and permitted by the supplying entity.”

8546 arrays of objects of this type are unlikely to be fully portable in any case. Solutions
8547 that treat them as magic cookies will be portable.

8548 {CHILD_MAX} precludes the possibility of a “toy implementation”, where there
8549 would only be one process.

8550 **ssize_t** This is intended to be a signed analog of **size_t**. The wording is such that an
8551 implementation may either choose to use a longer type or simply to use the signed
8552 version of the type that underlies **size_t**. All functions that return **ssize_t** (*read()*
8553 and *write()*) describe as “implementation-defined” the result of an input exceeding
8554 {SSIZE_MAX}. It is recognized that some implementations might have **ints** that
8555 are smaller than **size_t**. A conforming application would be constrained not to
8556 perform I/O in pieces larger than {SSIZE_MAX}, but a conforming application
8557 using extensions would be able to use the full range if the implementation
8558 provided an extended range, while still having a single type-compatible interface.

8559 The symbols **size_t** and **ssize_t** are also required in `<unistd.h>` to minimize the
8560 changes needed for calls to *read()* and *write()*. Implementors are reminded that it
8561 must be possible to include both `<sys/types.h>` and `<unistd.h>` in the same
8562 program (in either order) without error.

8563 **uid_t** Before the addition of this type, the data types used to represent these values
8564 varied throughout early proposals. The `<sys/stat.h>` header defined these values as
8565 type **short**, the `<passwd.h>` file (now `<pwd.h>` and `<grp.h>`) used an **int**, and
8566 *getuid()* returned an **int**. In response to a strong objection to the inconsistent
8567 definitions, all the types were switched to **uid_t**.

8568 In practice, those historical implementations that use varying types of this sort can
8569 typedef **uid_t** to **short** with no serious consequences.

8570 The problem associated with this change concerns object compatibility after
8571 structure size changes. Since most implementations will define **uid_t** as a short, the
8572 only substantive change will be a reduction in the size of the **passwd** structure.
8573 Consequently, implementations with an overriding concern for object
8574 compatibility can pad the structure back to its current size. For that reason, this
8575 problem was not considered critical enough to warrant the addition of a separate
8576 type to POSIX.1.

8577 The types **uid_t** and **gid_t** are magic cookies. There is no {UID_MAX} defined by
8578 POSIX.1, and no structure imposed on **uid_t** and **gid_t** other than that they be
8579 positive arithmetic types. (In fact, they could be **unsigned char**.) There is no
8580 maximum or minimum specified for the number of distinct user or group IDs.

8581 **B.3 System Interfaces**

8582 See the RATIONALE sections on the individual reference pages.

8583 **B.3.1 Examples for Spawn**8584 The following long examples are provided in the Rationale (Informative) volume of
8585 IEEE Std 1003.1-2001 as a supplement to the reference page for *posix_spawn()*.8586 **Example Library Implementation of Spawn**8587 The *posix_spawn()* or *posix_spawnnp()* functions provide the following:

- 8588 • Simply start a process executing a process image. This is the simplest application for process
8589 creation, and it may cover most executions of *fork()*.
- 8590 • Support I/O redirection, including pipes.
- 8591 • Run the child under a user and group ID in the domain of the parent.
- 8592 • Run the child at any priority in the domain of the parent.

8593 The *posix_spawn()* or *posix_spawnnp()* functions do not cover every possible use of the *fork()*
8594 function, but they do span the common applications: typical use by a shell and a login utility.8595 The price for an application is that before it calls *posix_spawn()* or *posix_spawnnp()*, the parent
8596 must adjust to a state that *posix_spawn()* or *posix_spawnnp()* can map to the desired state for the
8597 child. Environment changes require the parent to save some of its state and restore it afterwards.
8598 The effective behavior of a successful invocation of *posix_spawn()* is as if the operation were
8599 implemented with POSIX operations as follows:

```

8600 #include <sys/types.h>
8601 #include <stdlib.h>
8602 #include <stdio.h>
8603 #include <unistd.h>
8604 #include <sched.h>
8605 #include <fcntl.h>
8606 #include <signal.h>
8607 #include <errno.h>
8608 #include <string.h>
8609 #include <signal.h>

8610 /* #include <spawn.h> */
8611 /*****
8612  /* Things that could be defined in spawn.h */
8613  *****/
8614 typedef struct
8615 {
8616     short posix_attr_flags;
8617     #define POSIX_SPAWN_SETPGROUP          0x1
8618     #define POSIX_SPAWN_SETSIGMASK       0x2
8619     #define POSIX_SPAWN_SETSIGDEF       0x4
8620     #define POSIX_SPAWN_SETSCHEDULER    0x8
8621     #define POSIX_SPAWN_SETSCHEDPARAM   0x10
8622     #define POSIX_SPAWN_RESETIDS        0x20
8623     pid_t posix_attr_pgroup;
8624     sigset_t posix_attr_sigmask;
8625     sigset_t posix_attr_sigdefault;

```

```

8626         int posix_attr_schedpolicy;
8627         struct sched_param posix_attr_schedparam;
8628     }    posix_spawnattr_t;

8629     typedef char *posix_spawn_file_actions_t;

8630     int posix_spawn_file_actions_init(
8631         posix_spawn_file_actions_t *file_actions);
8632     int posix_spawn_file_actions_destroy(
8633         posix_spawn_file_actions_t *file_actions);
8634     int posix_spawn_file_actions_addclose(
8635         posix_spawn_file_actions_t *file_actions, int fildes);
8636     int posix_spawn_file_actions_adddup2(
8637         posix_spawn_file_actions_t *file_actions, int fildes,
8638         int newfildes);
8639     int posix_spawn_file_actions_addopen(
8640         posix_spawn_file_actions_t *file_actions, int fildes,
8641         const char *path, int oflag, mode_t mode);
8642     int posix_spawnattr_init(posix_spawnattr_t *attr);
8643     int posix_spawnattr_destroy(posix_spawnattr_t *attr);
8644     int posix_spawnattr_getflags(const posix_spawnattr_t *attr,
8645         short *lags);
8646     int posix_spawnattr_setflags(posix_spawnattr_t *attr, short flags);
8647     int posix_spawnattr_getpgroup(const posix_spawnattr_t *attr,
8648         pid_t *pgroup);
8649     int posix_spawnattr_setpgroup(posix_spawnattr_t *attr, pid_t pgroup);
8650     int posix_spawnattr_getschedpolicy(const posix_spawnattr_t *attr,
8651         int *schedpolicy);
8652     int posix_spawnattr_setschedpolicy(posix_spawnattr_t *attr,
8653         int schedpolicy);
8654     int posix_spawnattr_getschedparam(const posix_spawnattr_t *attr,
8655         struct sched_param *schedparam);
8656     int posix_spawnattr_setschedparam(posix_spawnattr_t *attr,
8657         const struct sched_param *schedparam);
8658     int posix_spawnattr_getsigmask(const posix_spawnattr_t *attr,
8659         sigset_t *sigmask);
8660     int posix_spawnattr_setsigmask(posix_spawnattr_t *attr,
8661         const sigset_t *sigmask);
8662     int posix_spawnattr_getdefault(const posix_spawnattr_t *attr,
8663         sigset_t *sigdefault);
8664     int posix_spawnattr_setsigdefault(posix_spawnattr_t *attr,
8665         const sigset_t *sigdefault);
8666     int posix_spawn(pid_t *pid, const char *path,
8667         const posix_spawn_file_actions_t *file_actions,
8668         const posix_spawnattr_t *attrp, char *const argv[],
8669         char *const envp[]);
8670     int posix_spawnvp(pid_t *pid, const char *file,
8671         const posix_spawn_file_actions_t *file_actions,
8672         const posix_spawnattr_t *attrp, char *const argv[],
8673         char *const envp[]);

8674     /*****
8675     /* Example posix_spawn() library routine */
8676     /*****

```

```

8677     int posix_spawn(pid_t *pid,
8678                   const char *path,
8679                   const posix_spawn_file_actions_t *file_actions,
8680                   const posix_spawnattr_t *attrp,
8681                   char *const argv[],
8682                   char *const envp[])
8683     {
8684         /* Create process */
8685         if ((*pid = fork()) == (pid_t) 0)
8686         {
8687             /* This is the child process */
8688             /* Worry about process group */
8689             if (attrp->posix_attr_flags & POSIX_SPAWN_SETPGROUP)
8690             {
8691                 /* Override inherited process group */
8692                 if (setpgid(0, attrp->posix_attr_pgroup) != 0)
8693                 {
8694                     /* Failed */
8695                     exit(127);
8696                 }
8697             }
8698             /* Worry about process signal mask */
8699             if (attrp->posix_attr_flags & POSIX_SPAWN_SETSIGMASK)
8700             {
8701                 /* Set the signal mask (can't fail) */
8702                 sigprocmask(SIG_SETMASK, &attrp->posix_attr_sigmask, NULL);
8703             }
8704             /* Worry about resetting effective user and group IDs */
8705             if (attrp->posix_attr_flags & POSIX_SPAWN_RESETIDS)
8706             {
8707                 /* None of these can fail for this case. */
8708                 setuid(getuid());
8709                 setgid(getgid());
8710             }
8711             /* Worry about defaulted signals */
8712             if (attrp->posix_attr_flags & POSIX_SPAWN_SETSIGDEF)
8713             {
8714                 struct sigaction deflt;
8715                 sigset_t all_signals;
8716
8717                 int s;
8718
8719                 /* Construct default signal action */
8720                 deflt.sa_handler = SIG_DFL;
8721                 deflt.sa_flags = 0;
8722
8723                 /* Construct the set of all signals */
8724                 sigfillset(&all_signals);
8725
8726                 /* Loop for all signals */
8727                 for (s = 0; sigismember(&all_signals, s); s++)
8728                 {
8729                     /* Signal to be defaulted? */

```

```

8726         if (sigismember(&attrp->posix_attr_sigdefault, s))
8727         {
8728             /* Yes; default this signal */
8729             if (sigaction(s, &deflt, NULL) == -1)
8730             {
8731                 /* Failed */
8732                 exit(127);
8733             }
8734         }
8735     }
8736 }

8737 /* Worry about the fds if they are to be mapped */
8738 if (file_actions != NULL)
8739 {
8740     /* Loop for all actions in object file_actions */
8741     /* (implementation dives beneath abstraction) */
8742     char *p = *file_actions;

8743     while (*p != '\0')
8744     {
8745         if (strncmp(p, "close(", 6) == 0)
8746         {
8747             int fd;

8748             if (sscanf(p + 6, "%d", &fd) != 1)
8749             {
8750                 exit(127);
8751             }
8752             if (close(fd) == -1)
8753                 exit(127);
8754         }
8755         else if (strncmp(p, "dup2(", 5) == 0)
8756         {
8757             int fd, newfd;

8758             if (sscanf(p + 5, "%d,%d", &fd, &newfd) != 2)
8759             {
8760                 exit(127);
8761             }
8762             if (dup2(fd, newfd) == -1)
8763                 exit(127);
8764         }
8765         else if (strncmp(p, "open(", 5) == 0)
8766         {
8767             int fd, oflag;
8768             mode_t mode;
8769             int tempfd;
8770             char path[1000];    /* Should be dynamic */
8771             char *q;

8772             if (sscanf(p + 5, "%d,", &fd) != 1)
8773             {
8774                 exit(127);
8775             }

```

```

8776         p = strchr(p, ',') + 1;
8777         q = strchr(p, '*');
8778         if (q == NULL)
8779             exit(127);
8780         strncpy(path, p, q - p);
8781         path[q - p] = '\0';
8782         if (sscanf(q + 1, "%o,%o", &oflag, &mode) != 2)
8783             {
8784                 exit(127);
8785             }
8786         if (close(fd) == -1)
8787             {
8788                 if (errno != EBADF)
8789                     exit(127);
8790             }
8791         tempfd = open(path, oflag, mode);
8792         if (tempfd == -1)
8793             exit(127);
8794         if (tempfd != fd)
8795             {
8796                 if (dup2(tempfd, fd) == -1)
8797                     {
8798                         exit(127);
8799                     }
8800                 if (close(tempfd) == -1)
8801                     {
8802                         exit(127);
8803                     }
8804             }
8805     }
8806     else
8807     {
8808         exit(127);
8809     }
8810     p = strchr(p, ')') + 1;
8811 }
8812
8813 /* Worry about setting new scheduling policy and parameters */
8814 if (attrp->posix_attr_flags & POSIX_SPAWN_SETSCHEDULER)
8815 {
8816     if (sched_setscheduler(0, attrp->posix_attr_schedpolicy,
8817         &attrp->posix_attr_schedparam) == -1)
8818     {
8819         exit(127);
8820     }
8821 }
8822 /* Worry about setting only new scheduling parameters */
8823 if (attrp->posix_attr_flags & POSIX_SPAWN_SETSCHEDPARAM)
8824 {
8825     if (sched_setparam(0, &attrp->posix_attr_schedparam) == -1)
8826     {

```

```

8827             exit(127);
8828         }
8829     }

8830     /* Now execute the program at path */
8831     /* Any fd that still has FD_CLOEXEC set will be closed */
8832     execve(path, argv, envp);
8833     exit(127);          /* exec failed */
8834 }
8835 else
8836 {
8837     /* This is the parent (calling) process */
8838     if (*pid == (pid_t) - 1)
8839         return errno;
8840     return 0;
8841 }
8842 }

8843 /*****
8844 /* Here is a crude but effective implementation of the */
8845 /* file action object operators which store actions as */
8846 /* concatenated token-separated strings.          */
8847 /*****
8848 /* Create object with no actions. */
8849 int posix_spawn_file_actions_init(
8850     posix_spawn_file_actions_t *file_actions)
8851 {
8852     *file_actions = malloc(sizeof(char));
8853     if (*file_actions == NULL)
8854         return ENOMEM;
8855     strcpy(*file_actions, "");
8856     return 0;
8857 }

8858 /* Free object storage and make invalid. */
8859 int posix_spawn_file_actions_destroy(
8860     posix_spawn_file_actions_t *file_actions)
8861 {
8862     free(*file_actions);
8863     *file_actions = NULL;
8864     return 0;
8865 }

8866 /* Add a new action string to object. */
8867 static int add_to_file_actions(
8868     posix_spawn_file_actions_t *file_actions, char *new_action)
8869 {
8870     *file_actions = realloc
8871     (*file_actions, strlen(*file_actions) + strlen(new_action) + 1);
8872     if (*file_actions == NULL)
8873         return ENOMEM;
8874     strcat(*file_actions, new_action);
8875     return 0;
8876 }

```

```

8877     /* Add a close action to object. */
8878     int posix_spawn_file_actions_addclose(
8879         posix_spawn_file_actions_t *file_actions, int fildes)
8880     {
8881         char temp[100];
8882
8883         sprintf(temp, "close(%d)", fildes);
8884         return add_to_file_actions(file_actions, temp);
8885     }
8886
8887     /* Add a dup2 action to object. */
8888     int posix_spawn_file_actions_adddup2(
8889         posix_spawn_file_actions_t *file_actions, int fildes,
8890         int newfildes)
8891     {
8892         char temp[100];
8893
8894         sprintf(temp, "dup2(%d,%d)", fildes, newfildes);
8895         return add_to_file_actions(file_actions, temp);
8896     }
8897
8898     /* Add an open action to object. */
8899     int posix_spawn_file_actions_addopen(
8900         posix_spawn_file_actions_t *file_actions, int fildes,
8901         const char *path, int oflag, mode_t mode)
8902     {
8903         char temp[100];
8904
8905         sprintf(temp, "open(%d,%s*%o,%o)", fildes, path, oflag, mode);
8906         return add_to_file_actions(file_actions, temp);
8907     }
8908
8909     /******
8910     /* Here is a crude but effective implementation of the */
8911     /* spawn attributes object functions which manipulate */
8912     /* the individual attributes. */
8913     /******
8914     /* Initialize object with default values. */
8915     int posix_spawnattr_init(posix_spawnattr_t *attr)
8916     {
8917         attr->posix_attr_flags = 0;
8918         attr->posix_attr_pgroup = 0;
8919         /* Default value of signal mask is the parent's signal mask; */
8920         /* other values are also allowed */
8921         sigprocmask(0, NULL, &attr->posix_attr_sigmask);
8922         sigemptyset(&attr->posix_attr_sigdefault);
8923         /* Default values of scheduling attr inherited from the parent; */
8924         /* other values are also allowed */
8925         attr->posix_attr_schedpolicy = sched_getscheduler(0);
8926         sched_getparam(0, &attr->posix_attr_schedparam);
8927         return 0;
8928     }
8929
8930     int posix_spawnattr_destroy(posix_spawnattr_t *attr)
8931     {
8932         /* No action needed */

```

```
8926     return 0;
8927 }
8928 int posix_spawnattr_getflags(const posix_spawnattr_t *attr,
8929     short *flags)
8930 {
8931     *flags = attr->posix_attr_flags;
8932     return 0;
8933 }
8934 int posix_spawnattr_setflags(posix_spawnattr_t *attr, short flags)
8935 {
8936     attr->posix_attr_flags = flags;
8937     return 0;
8938 }
8939 int posix_spawnattr_getpgroup(const posix_spawnattr_t *attr,
8940     pid_t *pgroup)
8941 {
8942     *pgroup = attr->posix_attr_pgroup;
8943     return 0;
8944 }
8945 int posix_spawnattr_setpgroup(posix_spawnattr_t *attr, pid_t pgroup)
8946 {
8947     attr->posix_attr_pgroup = pgroup;
8948     return 0;
8949 }
8950 int posix_spawnattr_getschedpolicy(const posix_spawnattr_t *attr,
8951     int *schedpolicy)
8952 {
8953     *schedpolicy = attr->posix_attr_schedpolicy;
8954     return 0;
8955 }
8956 int posix_spawnattr_setschedpolicy(posix_spawnattr_t *attr,
8957     int schedpolicy)
8958 {
8959     attr->posix_attr_schedpolicy = schedpolicy;
8960     return 0;
8961 }
8962 int posix_spawnattr_getschedparam(const posix_spawnattr_t *attr,
8963     struct sched_param *schedparam)
8964 {
8965     *schedparam = attr->posix_attr_schedparam;
8966     return 0;
8967 }
8968 int posix_spawnattr_setschedparam(posix_spawnattr_t *attr,
8969     const struct sched_param *schedparam)
8970 {
8971     attr->posix_attr_schedparam = *schedparam;
8972     return 0;
8973 }
```



```

8974     int posix_spawnattr_getsigmask(const posix_spawnattr_t *attr,
8975         sigset_t *sigmask)
8976     {
8977         *sigmask = attr->posix_attr_sigmask;
8978         return 0;
8979     }

8980     int posix_spawnattr_setsigmask(posix_spawnattr_t *attr,
8981         const sigset_t *sigmask)
8982     {
8983         attr->posix_attr_sigmask = *sigmask;
8984         return 0;
8985     }

8986     int posix_spawnattr_getsigdefault(const posix_spawnattr_t *attr,
8987         sigset_t *sigdefault)
8988     {
8989         *sigdefault = attr->posix_attr_sigdefault;
8990         return 0;
8991     }

8992     int posix_spawnattr_setsigdefault(posix_spawnattr_t *attr,
8993         const sigset_t *sigdefault)
8994     {
8995         attr->posix_attr_sigdefault = *sigdefault;
8996         return 0;
8997     }

```

8998 **I/O Redirection with Spawn**

8999 I/O redirection with *posix_spawn()* or *posix_spawnp()* is accomplished by crafting a *file_actions*
9000 argument to effect the desired redirection. Such a redirection follows the general outline of the
9001 following example:

```

9002     /* To redirect new standard output (fd 1) to a file, */
9003     /* and redirect new standard input (fd 0) from my fd socket_pair[1], */
9004     /* and close my fd socket_pair[0] in the new process. */
9005     posix_spawn_file_actions_t file_actions;
9006     posix_spawn_file_actions_init(&file_actions);
9007     posix_spawn_file_actions_addopen(&file_actions, 1, "newout", ...);
9008     posix_spawn_file_actions_dup2(&file_actions, socket_pair[1], 0);
9009     posix_spawn_file_actions_close(&file_actions, socket_pair[0]);
9010     posix_spawn_file_actions_close(&file_actions, socket_pair[1]);
9011     posix_spawn(..., &file_actions, ...);
9012     posix_spawn_file_actions_destroy(&file_actions);

```

9013 Spawning a Process Under a New User ID

9014 Spawning a process under a new user ID follows the outline shown in the following example:

```
9015            Save = getuid();  
9016            setuid(newid);  
9017            posix_spawn(...);  
9018            setuid(Save);
```

9019 / *Rationale (Informative)*

9020 **Part C:**

9021 **Shell and Utilities**

9022 *The Open Group*

9023 *The Institute of Electrical and Electronics Engineers, Inc.*

Rationale for Shell and Utilities

9024

9025 **C.1 Introduction**9026 **C.1.1 Scope**

9027 Refer to Section A.1.1 (on page 3).

9028 **C.1.2 Conformance**

9029 Refer to Section A.2 (on page 9).

9030 **C.1.3 Normative References**

9031 There is no additional rationale provided for this section.

9032 **C.1.4 Change History**9033 The change history is provided as an informative section, to track changes from previous issues
9034 of IEEE Std 1003.1-2001.9035 The following sections describe changes made to the Shell and Utilities volume of
9036 IEEE Std 1003.1-2001 since Issue 5 of the base document. The CHANGE HISTORY section for
9037 each utility describes technical changes made to that utility from Issue 5. Changes between
9038 earlier issues of the base document and Issue 5 are not included.9039 The change history between Issue 5 and Issue 6 also lists the changes since the
9040 ISO POSIX-2: 1993 standard.9041 **Changes from Issue 5 to Issue 6 (IEEE Std 1003.1-2001)**9042 The following list summarizes the major changes that were made in the Shell and Utilities
9043 volume of IEEE Std 1003.1-2001 from Issue 5 to Issue 6:9044 • This volume of IEEE Std 1003.1-2001 is extensively revised so that it can be both an IEEE
9045 POSIX Standard and an Open Group Technical Standard.

9046 • The terminology has been reworked to meet the style requirements.

9047 • Shading notation and margin codes are introduced for identification of options within the
9048 volume.9049 • This volume of IEEE Std 1003.1-2001 is updated to mandate support of FIPS 151-2. The
9050 following changes were made:9051 — Support is mandated for the capabilities associated with the following symbolic
9052 constants:9053 `_POSIX_CHOWN_RESTRICTED`9054 `_POSIX_JOB_CONTROL`9055 `_POSIX_SAVED_IDS`9056 — In the environment for the login shell, the environment variables *LOGNAME* and *HOME*
9057 shall be defined and have the properties described in the Base Definitions volume of

- 9058 IEEE Std 1003.1-2001, Chapter 7, Locale.
- 9059 • This volume of IEEE Std 1003.1-2001 is updated to align with some features of the Single
 - 9060 UNIX Specification.
 - 9061 • A new section on Utility Limits is added.
 - 9062 • A section on the Relationships to Other Documents is added.
 - 9063 • Concepts and definitions have been moved to a separate volume.
 - 9064 • A RATIONALE section is added to each reference page.
 - 9065 • The *c99* utility is added as a replacement for *c89*, which is withdrawn in this issue.
 - 9066 • IEEE Std 1003.2d-1994 is incorporated, adding the *qalter*, *qdel*, *qhold*, *qmove*, *qmsg*, *qrerun*, *qrls*,
 - 9067 *qselect*, *qsig*, *qstat*, and *qsub* utilities.
 - 9068 • IEEE P1003.2b draft standard is incorporated, making extensive updates and adding the *iconv*
 - 9069 utility.
 - 9070 • IEEE PASC Interpretations are applied.
 - 9071 • The Open Group's corrigenda and resolutions are applied.

9072 New Features in Issue 6

9073 The following table lists the new utilities introduced since the ISO POSIX-2:1993 standard (as
 9074 modified by IEEE Std 1003.2d-1994). Apart from the *c99* and *iconv* utilities, these are all part of
 9075 the XSI extension.

9076
 9077

New Utilities in Issue 6							
<i>admin</i>	<i>compress</i>	<i>gencat</i>	<i>ipcrm</i>	<i>nl</i>	<i>tsort</i>	<i>unlink</i>	<i>val</i>
<i>c99</i>	<i>cxref</i>	<i>get</i>	<i>ipcs</i>	<i>prs</i>	<i>ulimit</i>	<i>uucp</i>	<i>what</i>
<i>cal</i>	<i>delta</i>	<i>hash</i>	<i>link</i>	<i>sact</i>	<i>uncompress</i>	<i>uustat</i>	<i>zcat</i>
<i>cflow</i>	<i>fuser</i>	<i>iconv</i>	<i>m4</i>	<i>sccs</i>	<i>unget</i>	<i>uux</i>	

9082 C.1.5 Terminology

9083 Refer to Section A.1.4 (on page 5).

9084 C.1.6 Definitions

9085 Refer to Section A.3 (on page 13).

9086 C.1.7 Relationship to Other Documents

9087 C.1.7.1 System Interfaces

9088 It has been pointed out that the Shell and Utilities volume of IEEE Std 1003.1-2001 assumes that
 9089 a great deal of functionality from the System Interfaces volume of IEEE Std 1003.1-2001 is
 9090 present, but never states exactly how much (and strictly does not need to since both are
 9091 mandated on a conforming system). This section is an attempt to clarify the assumptions.

9092 File Removal

9093 This is intended to be a summary of the *unlink()* and *rmdir()* requirements. Note that it is
9094 possible using the *unlink()* function for item 4. to occur.

9095 C.1.7.2 Concepts Derived from the ISO C Standard

9096 This section was introduced to address the issue that there was insufficient detail presented by
9097 such utilities as *awk* or *sh* about their procedural control statements and their methods of
9098 performing arithmetic functions.

9099 The ISO C standard was selected as a model because most historical implementations of the
9100 standard utilities were written in C. Thus, it was more likely that they would act in the desired
9101 manner without modification.

9102 Using the ISO C standard is primarily a notational convenience so that the many procedural
9103 languages in the Shell and Utilities volume of IEEE Std 1003.1-2001 would not have to be
9104 rigorously described in every aspect. Its selection does not require that the standard utilities be
9105 written in Standard C; they could be written in Common Usage C, Ada, Pascal, assembler
9106 language, or anything else.

9107 The sizes of the various numeric values refer to C-language data types that are allowed to be
9108 different sizes by the ISO C standard. Thus, like a C-language application, a shell application
9109 cannot rely on their exact size. However, it can rely on their minimum sizes expressed in the
9110 ISO C standard, such as {LONG_MAX} for a **long** type.

9111 The behavior on overflow is undefined for ISO C standard arithmetic. Therefore, the standard
9112 utilities can use “bignum” representation for integers so that there is no fixed maximum unless
9113 otherwise stated in the utility description. Similarly, standard utilities can use infinite-precision
9114 representations for floating-point arithmetic, as long as these representations exceed the ISO C
9115 standard requirements.

9116 This section addresses only the issue of semantics; it is not intended to specify syntax. For
9117 example, the ISO C standard requires that 0L be recognized as an integer constant equal to zero,
9118 but utilities such as *awk* and *sh* are not required to recognize 0L (though they are allowed to, as
9119 an extension).

9120 The ISO C standard requires that a C compiler must issue a diagnostic for constants that are too
9121 large to represent. Most standard utilities are not required to issue these diagnostics; for
9122 example, the command:

```
9123 diff -C 2147483648 file1 file2
```

9124 has undefined behavior, and the *diff* utility is not required to issue a diagnostic even if the
9125 number 2 147 483 648 cannot be represented.

9126 C.1.8 Portability

9127 Refer to Section A.1.5 (on page 8).

9128 C.1.8.1 Codes

9129 Refer to Section A.1.5.1 (on page 8).

9130 **C.1.9 Utility Limits**

9131 This section grew out of an idea that originated with the original POSIX.1, in the tables of system
 9132 limits for the *sysconf()* and *pathconf()* functions. The idea being that a conforming application
 9133 can be written to use the most restrictive values that a minimal system can provide, but it should
 9134 not have to. The values provided represent compromises so that some vendors can use
 9135 historically limited versions of UNIX system utilities. They are the highest values that a strictly
 9136 conforming application can assume, given no other information.

9137 However, by using the *getconf* utility or the *sysconf()* function, the elegant application can be
 9138 tailored to more liberal values on some of the specific instances of specific implementations.

9139 There is no explicitly stated requirement that an implementation provide finite limits for any of
 9140 these numeric values; the implementation is free to provide essentially unbounded capabilities
 9141 (where it makes sense), stopping only at reasonable points such as {ULONG_MAX} (from the
 9142 ISO C standard). Therefore, applications desiring to tailor themselves to the values on a
 9143 particular implementation need to be ready for possibly huge values; it may not be a good idea
 9144 to allocate blindly a buffer for an input line based on the value of {LINE_MAX}, for instance.
 9145 However, unlike the System Interfaces volume of IEEE Std 1003.1-2001, there is no set of limits
 9146 that return a special indication meaning “unbounded”. The implementation should always
 9147 return an actual number, even if the number is very large.

9148 The statement:

9149 “It is not guaranteed that the application ...”

9150 is an indication that many of these limits are designed to ensure that implementors design their
 9151 utilities without arbitrary constraints related to unimaginative programming. There are certainly
 9152 conditions under which combinations of options can cause failures that would not render an
 9153 implementation non-conforming. For example, {EXPR_NEST_MAX} and {ARG_MAX} could
 9154 collide when expressions are large; combinations of {BC_SCALE_MAX} and {BC_DIM_MAX}
 9155 could exceed virtual memory.

9156 In the Shell and Utilities volume of IEEE Std 1003.1-2001, the notion of a limit being guaranteed
 9157 for the process lifetime, as it is in the System Interfaces volume of IEEE Std 1003.1-2001, is not as
 9158 useful to a shell script. The *getconf* utility is probably a process itself, so the guarantee would be
 9159 without value. Therefore, the Shell and Utilities volume of IEEE Std 1003.1-2001 requires the
 9160 guarantee to be for the session lifetime. This will mean that many vendors will either return very
 9161 conservative values or possibly implement *getconf* as a built-in.

9162 It may seem confusing to have limits that apply only to a single utility grouped into one global
 9163 section. However, the alternative, which would be to disperse them out into their utility
 9164 description sections, would cause great difficulty when *sysconf()* and *getconf* were described.
 9165 Therefore, the standard developers chose the global approach.

9166 Each language binding could provide symbol names that are slightly different from those shown
 9167 here. For example, the C-Language Binding option adds a leading underscore to the symbols as a
 9168 prefix.

9169 The following comments describe selection criteria for the symbols and their values:

9170 {ARG_MAX}

9171 This is defined by the System Interfaces volume of IEEE Std 1003.1-2001. Unfortunately, it is
 9172 very difficult for a conforming application to deal with this value, as it does not know how
 9173 much of its argument space is being consumed by the environment variables of the user.

9174 {BC_BASE_MAX}
 9175 {BC_DIM_MAX}
 9176 {BC_SCALE_MAX}
 9177 These were originally one value, {BC_SCALE_MAX}, but it was unreasonable to link all
 9178 three concepts into one limit.

9179 {CHILD_MAX}
 9180 This is defined by the System Interfaces volume of IEEE Std 1003.1-2001.

9181 {COLL_WEIGHTS_MAX}
 9182 The weights assigned to **order** can be considered as “passes” through the collation
 9183 algorithm.

9184 {EXPR_NEST_MAX}
 9185 The value for expression nesting was borrowed from the ISO C standard.

9186 {LINE_MAX}
 9187 This is a global limit that affects all utilities, unless otherwise noted. The {MAX_CANON}
 9188 value from the System Interfaces volume of IEEE Std 1003.1-2001 may further limit input
 9189 lines from terminals. The {LINE_MAX} value was the subject of much debate and is a
 9190 compromise between those who wished to have unlimited lines and those who understood
 9191 that many historical utilities were written with fixed buffers. Frequently, utility writers
 9192 selected the UNIX system constant BUFSIZ to allocate these buffers; therefore, some utilities
 9193 were limited to 512 bytes for I/O lines, while others achieved 4 096 bytes or greater.

9194 It should be noted that {LINE_MAX} applies only to input line length; there is no
 9195 requirement in IEEE Std 1003.1-2001 that limits the length of output lines. Utilities such as
 9196 *awk*, *sed*, and *paste* could theoretically construct lines longer than any of the input lines they
 9197 received, depending on the options used or the instructions from the application. They are
 9198 not required to truncate their output to {LINE_MAX}. It is the responsibility of the
 9199 application to deal with this. If the output of one of those utilities is to be piped into another
 9200 of the standard utilities, line length restrictions will have to be considered; the *fold* utility,
 9201 among others, could be used to ensure that only reasonable line lengths reach utilities or
 9202 applications.

9203 {LINK_MAX}
 9204 This is defined by the System Interfaces volume of IEEE Std 1003.1-2001.

9205 {MAX_CANON}
 9206 {MAX_INPUT}
 9207 {NAME_MAX}
 9208 {NGROUPS_MAX}
 9209 {OPEN_MAX}
 9210 {PATH_MAX}
 9211 {PIPE_BUF}

9212 These limits are defined by the System Interfaces volume of IEEE Std 1003.1-2001. Note that
 9213 the byte lengths described by some of these values continue to represent bytes, even if the
 9214 applicable character set uses a multi-byte encoding.

9215 {RE_DUP_MAX}
 9216 The value selected is consistent with historical practice. Although the name implies that it
 9217 applies to all REs, only BREs use the interval notation $\{m,n\}$ addressed by this limit.

9218 {POSIX2_SYMLINKS}
 9219 The {POSIX2_SYMLINKS} variable indicates that the underlying operating system supports
 9220 the creation of symbolic links in specific directories. Many of the utilities defined in
 9221 IEEE Std 1003.1-2001 that deal with symbolic links do not depend on this value. For

9222 example, a utility that follows symbolic links (or does not, as the case may be) will only be
 9223 affected by a symbolic link if it encounters one. Presumably, a file system that does not
 9224 support symbolic links will not contain any. This variable does affect such utilities as *ln -s*
 9225 and *pax* that attempt to create symbolic links.

9226 {POSIX2_SYMLINKS} was developed even though there is no comparable configuration
 9227 value for the system interfaces.

9228 There are different limits associated with command lines and input to utilities, depending on the
 9229 method of invocation. In the case of a C program *exec*-ing a utility, {ARG_MAX} is the
 9230 underlying limit. In the case of the shell reading a script and *exec*-ing a utility, {LINE_MAX}
 9231 limits the length of lines the shell is required to process, and {ARG_MAX} will still be a limit. If a
 9232 user is entering a command on a terminal to the shell, requesting that it invoke the utility,
 9233 {MAX_INPUT} may restrict the length of the line that can be given to the shell to a value below
 9234 {LINE_MAX}.

9235 When an option is supported, *getconf* returns a value of 1. For example, when C development is
 9236 supported:

```
9237     if [ "$(getconf POSIX2_C_DEV)" -eq 1 ]; then
9238         echo C supported
9239     fi
```

9240 The *sysconf()* function in the C-Language Binding option would return 1.

9241 The following comments describe selection criteria for the symbols and their values:

```
9242 POSIX2_C_BIND
9243 POSIX2_C_DEV
9244 POSIX2_FORT_DEV
9245 POSIX2_FORT_RUN
9246 POSIX2_SW_DEV
9247 POSIX2_UPE
```

9248 It is possible for some (usually privileged) operations to remove utilities that support these
 9249 options or otherwise to render these options unsupported. The header files, the *sysconf()*
 9250 function, or the *getconf* utility will not necessarily detect such actions, in which case they
 9251 should not be considered as rendering the implementation non-conforming. A test suite
 9252 should not attempt tests such as:

```
9253     rm /usr/bin/c99
9254     getconf POSIX2_C_DEV
```

```
9255 POSIX2_LOCALEDEF
```

9256 This symbol was introduced to allow implementations to restrict supported locales to only
 9257 those supplied by the implementation.

9258 C.1.10 Grammar Conventions

9259 There is no additional rationale provided for this section.

9260 **C.1.11 Utility Description Defaults**

9261 This section is arranged with headings in the same order as all the utility descriptions. It is a
9262 collection of related and unrelated information concerning:

- 9263 1. The default actions of utilities
- 9264 2. The meanings of notations used in IEEE Std 1003.1-2001 that are specific to individual
9265 utility sections

9266 Although this material may seem out of place here, it is important that this information appear
9267 before any of the utilities to be described later.

9268 **NAME**

9269 There is no additional rationale provided for this section.

9270 **SYNOPSIS**

9271 There is no additional rationale provided for this section.

9272 **DESCRIPTION**

9273 There is no additional rationale provided for this section.

9274 **OPTIONS**

9275 Although it has not always been possible, the standard developers tried to avoid repeating
9276 information to reduce the risk that duplicate explanations could each be modified differently.

9277 The need to recognize `--` is required because conforming applications need to shield their
9278 operands from any arbitrary options that the implementation may provide as an extension. For
9279 example, if the standard utility *foo* is listed as taking no options, and the application needed to
9280 give it a pathname with a leading hyphen, it could safely do it as:

9281 `foo -- -myfile`

9282 and avoid any problems with `-m` used as an extension.

9283 **OPERANDS**

9284 The usage of `-` is never shown in the SYNOPSIS. Similarly, the usage of `--` is never shown.

9285 The requirement for processing operands in command-line order is to avoid a “WeirdNIX”
9286 utility that might choose to sort the input files alphabetically, by size, or by directory order.
9287 Although this might be acceptable for some utilities, in general the programmer has a right to
9288 know exactly what order will be chosen.

9289 Some of the standard utilities take multiple *file* operands and act as if they were processing the
9290 concatenation of those files. For example:

9291 `asa file1 file2`

9292 and:

9293 `cat file1 file2 | asa`

9294 have similar results when questions of file access, errors, and performance are ignored. Other
9295 utilities such as *grep* or *wc* have completely different results in these two cases. This latter type of
9296 utility is always identified in its DESCRIPTION or OPERANDS sections, whereas the former is
9297 not. Although it might be possible to create a general assertion about the former case, the

9298 following points must be addressed:

- 9299
- Access times for the files might be different in the operand case *versus* the *cat* case.
 - The utility may have error messages that are cognizant of the input filename, and this added value should not be suppressed. (As an example, *awk* sets a variable with the filename at each file boundary.)
- 9300
- 9301
- 9302

9303 **STDIN**

9304 There is no additional rationale provided for this section.

9305 **INPUT FILES**

9306 A conforming application cannot assume the following three commands are equivalent:

```
9307 tail -n +2 file
9308 (sed -n 1q; cat) < file
9309 cat file | (sed -n 1q; cat)
```

9310 The second command is equivalent to the first only when the file is seekable. In the third
 9311 command, if the file offset in the open file description were not unspecified, *sed* would have to be
 9312 implemented so that it read from the pipe 1 byte at a time or it would have to employ some
 9313 method to seek backwards on the pipe. Such functionality is not defined currently in POSIX.1
 9314 and does not exist on all historical systems. Other utilities, such as *head*, *read*, and *sh*, have similar
 9315 properties, so the restriction is described globally in this section.

9316 The definition of “text file” is strictly enforced for input to the standard utilities; very few of
 9317 them list exceptions to the undefined results called for here. (Of course, “undefined” here does
 9318 not mean that historical implementations necessarily have to change to start indicating error
 9319 conditions. Conforming applications cannot rely on implementations succeeding or failing when
 9320 non-text files are used.)

9321 The utilities that allow line continuation are generally those that accept input languages, rather
 9322 than pure data. It would be unusual for an input line of this type to exceed {LINE_MAX} bytes
 9323 and unreasonable to require that the implementation allow unlimited accumulation of multiple
 9324 lines, each of which could reach {LINE_MAX}. Thus, for a conforming application the total of all
 9325 the continued lines in a set cannot exceed {LINE_MAX}.

9326 The format description is intended to be sufficiently rigorous to allow other applications to
 9327 generate these input files. However, since <blank>s can legitimately be included in some of the
 9328 fields described by the standard utilities, particularly in locales other than the POSIX locale, this
 9329 intent is not always realized.

9330 **ENVIRONMENT VARIABLES**

9331 There is no additional rationale provided for this section.

9332 **ASYNCHRONOUS EVENTS**

9333 Because there is no language prohibiting it, a utility is permitted to catch a signal, perform some
 9334 additional processing (such as deleting temporary files), restore the default signal action (or
 9335 action inherited from the parent process), and resignal itself.

9336 **STDOUT**

9337 The format description is intended to be sufficiently rigorous to allow post-processing of output
9338 by other programs, particularly by an *awk* or *lex* parser.

9339 **STDERR**

9340 This section does not describe error messages that refer to incorrect operation of the utility.
9341 Consider a utility that processes program source code as its input. This section is used to
9342 describe messages produced by a correctly operating utility that encounters an error in the
9343 program source code on which it is processing. However, a message indicating that the utility
9344 had insufficient memory in which to operate would not be described.

9345 Some utilities have traditionally produced warning messages without returning a non-zero exit
9346 status; these are specifically noted in their sections. Other utilities shall not write to standard
9347 error if they complete successfully, unless the implementation provides some sort of extension
9348 to increase the verbosity or debugging level.

9349 The format descriptions are intended to be sufficiently rigorous to allow post-processing of
9350 output by other programs.

9351 **OUTPUT FILES**

9352 The format description is intended to be sufficiently rigorous to allow post-processing of output
9353 by other programs, particularly by an *awk* or *lex* parser.

9354 Receipt of the SIGQUIT signal should generally cause termination (unless in some debugging
9355 mode) that would bypass any attempted recovery actions.

9356 **EXTENDED DESCRIPTION**

9357 There is no additional rationale provided for this section.

9358 **EXIT STATUS**

9359 Note the additional discussion of exit values in *Exit Status for Commands* in the *sh* utility. It
9360 describes requirements for returning exit values greater than 125.

9361 A utility may list zero as a successful return, 1 as a failure for a specific reason, and greater than
9362 1 as “an error occurred”. In this case, unspecified conditions may cause a 2 or 3, or other value,
9363 to be returned. A strictly conforming application should be written so that it tests for successful
9364 exit status values (zero in this case), rather than relying upon the single specific error value listed
9365 in IEEE Std 1003.1-2001. In that way, it will have maximum portability, even on implementations
9366 with extensions.

9367 The standard developers are aware that the general non-enumeration of errors makes it difficult
9368 to write test suites that test the *incorrect* operation of utilities. There are some historical
9369 implementations that have expended effort to provide detailed status messages and a helpful
9370 environment to bypass or explain errors, such as prompting, retrying, or ignoring unimportant
9371 syntax errors; other implementations have not. Since there is no realistic way to mandate system
9372 behavior in cases of undefined application actions or system problems—in a manner acceptable
9373 to all cultures and environments—attention has been limited to the correct operation of utilities
9374 by the conforming application. Furthermore, the conforming application does not need detailed
9375 information concerning errors that it caused through incorrect usage or that it cannot correct.

9376 There is no description of defaults for this section because all of the standard utilities specify
9377 something (or explicitly state “Unspecified”) for exit status.

9378 **CONSEQUENCES OF ERRORS**

9379 Several actions are possible when a utility encounters an error condition, depending on the
 9380 severity of the error and the state of the utility. Included in the possible actions of various
 9381 utilities are: deletion of temporary or intermediate work files; deletion of incomplete files; and
 9382 validity checking of the file system or directory.

9383 The text about recursive traversing is meant to ensure that utilities such as *find* process as many
 9384 files in the hierarchy as they can. They should not abandon all of the hierarchy at the first error
 9385 and resume with the next command-line operand, but should attempt to keep going.

9386 **APPLICATION USAGE**

9387 This section provides additional caveats, issues, and recommendations to the developer.

9388 **EXAMPLES**

9389 This section provides sample usage.

9390 **RATIONALE**

9391 There is no additional rationale provided for this section.

9392 **FUTURE DIRECTIONS**

9393 FUTURE DIRECTIONS sections act as pointers to related work that may impact the interface in
 9394 the future, and often cautions the developer to architect the code to account for a change in this
 9395 area. Note that a future directions statement should not be taken as a commitment to adopt a
 9396 feature or interface in the future.

9397 **SEE ALSO**

9398 There is no additional rationale provided for this section.

9399 **CHANGE HISTORY**

9400 There is no additional rationale provided for this section.

9401 **C.1.12 Considerations for Utilities in Support of Files of Arbitrary Size**

9402 This section is intended to clarify the requirements for utilities in support of large files.

9403 The utilities listed in this section are utilities which are used to perform administrative tasks
 9404 such as to create, move, copy, remove, change the permissions, or measure the resources of a
 9405 file. They are useful both as end-user tools and as utilities invoked by applications during
 9406 software installation and operation.

9407 The *chgrp*, *chmod*, *chown*, *ln*, and *rm* utilities probably require use of large file-capable versions of
 9408 *stat()*, *lstat()*, *ftw()*, and the **stat** structure.

9409 The *cat*, *cksum*, *cmp*, *cp*, *dd*, *mv*, *sum*, and *touch* utilities probably require use of large file-capable
 9410 versions of *creat()*, *open()*, and *fopen()*.

9411 The *cat*, *cksum*, *cmp*, *dd*, *df*, *du*, *ls*, and *sum* utilities may require writing large integer values. For
 9412 example:

- 9413 • The *cat* utility might have a *-n* option which counts <newline>s.
- 9414 • The *cksum* and *ls* utilities report file sizes.

9415 • The *cmp* utility reports the line number at which the first difference occurs, and also has a *-l*
9416 option which reports file offsets.

9417 • The *dd*, *df*, *du*, *ls*, and *sum* utilities report block counts.

9418 The *dd*, *find*, and *test* utilities may need to interpret command arguments that contain 64-bit
9419 values. For *dd*, the arguments include *skip=n*, *seek=n*, and *count=n*. For *find*, the arguments
9420 include *-sizen*. For *test*, the arguments are those associated with algebraic comparisons.

9421 The *df* utility might need to access large file systems with *statvfs()*.

9422 The *ulimit* utility will need to use large file-capable versions of *getrlimit()* and *setrlimit()* and be
9423 able to read and write large integer values.

9424 C.1.13 Built-In Utilities

9425 All of these utilities can be *exec*-ed. There is no requirement that these utilities are actually built
9426 into the shell itself, but many shells need the capability to do so because the Shell and Utilities
9427 volume of IEEE Std 1003.1-2001, Section 2.9.1.1, Command Search and Execution requires that
9428 they be found prior to the *PATH* search. The shell could satisfy its requirements by keeping a list
9429 of the names and directly accessing the file-system versions regardless of *PATH*. Providing all of
9430 the required functionality for those such as *cd* or *read* would be more difficult.

9431 There were originally three justifications for allowing the omission of *exec*-able versions:

9432 1. It would require wasting space in the file system, at the expense of very small systems.
9433 However, it has been pointed out that all 16 utilities in the table can be provided with 16
9434 links to a single-line shell script:

9435 `$0 "$@"`

9436 2. It is not logical to require invocation of utilities such as *cd* because they have no value
9437 outside the shell environment or cannot be useful in a child process. However, counter-
9438 examples always seemed to be available for even the most unusual cases:

9439 `find . -type d -exec cd {} \; -exec foo {} \;`
9440 (which invokes “foo” on accessible directories)

9441 `ps ... | sed ... | xargs kill`

9442 `find . -exec true \; -a ...`
9443 (where “true” is used for temporary debugging)

9444 3. It is confusing to have a utility such as *kill* that can easily be in the file system in the base
9445 standard, but that requires built-in status for the User Portability Utilities option (for the %
9446 job control job ID notation). It was decided that it was more appropriate to describe the
9447 required functionality (rather than the implementation) to the system implementors and
9448 let them decide how to satisfy it.

9449 On the other hand, it was realized that any distinction like this between utilities was not useful
9450 to applications, and that the cost to correct it was small. These arguments were ultimately the
9451 most effective.

9452 There were varying reasons for including utilities in the table of built-ins:

9453 *alias*, *fc*, *unalias*

9454 The functionality of these utilities is performed more simply within the shell itself and that
9455 is the model most historical implementations have used.

9456 *bg*, *fg*, *jobs*

9457 All of the job control-related utilities are eligible for built-in status because that is the model

9458 most historical implementations have used.

9459 *cd, getopts, newgrp, read, umask, wait*

9460 The functionality of these utilities is performed more simply within the context of the
 9461 current process. An example can be taken from the usage of the *cd* utility. The purpose of
 9462 the *cd* utility is to change the working directory for subsequent operations. The actions of *cd*
 9463 affect the process in which *cd* is executed and all subsequent child processes of that process.
 9464 Based on the POSIX standard process model, changes in the process environment of a child
 9465 process have no effect on the parent process. If the *cd* utility were executed from a child
 9466 process, the working directory change would be effective only in the child process. Child
 9467 processes initiated subsequent to the child process that executed the *cd* utility would not
 9468 have a changed working directory relative to the parent process.

9469 *command*

9470 This utility was placed in the table primarily to protect scripts that are concerned about
 9471 their *PATH* being manipulated. The “secure” shell script example in the *command* utility in
 9472 the Shell and Utilities volume of IEEE Std 1003.1-2001 would not be possible if a *PATH*
 9473 change retrieved an alien version of *command*. (An alternative would have been to
 9474 implement *getconf* as a built-in, but the standard developers considered that it carried too
 9475 many changing configuration strings to require in the shell.)

9476 *kill* Since *kill* provides optional job control functionality using shell notation (%1, %2, and so on),
 9477 some implementations would find it extremely difficult to provide this outside the shell.

9478 *true, false*

9479 These are in the table as a courtesy to programmers who wish to use the “while true”
 9480 shell construct without protecting *true* from *PATH* searches. (It is acknowledged that
 9481 “while :” also works, but the idiom with *true* is historically pervasive.)

9482 All utilities, including those in the table, are accessible via the *system()* and *popen()* functions in
 9483 the System Interfaces volume of IEEE Std 1003.1-2001. There are situations where the return
 9484 functionality of *system()* and *popen()* is not desirable. Applications that require the exit status of
 9485 the invoked utility will not be able to use *system()* or *popen()*, since the exit status returned is
 9486 that of the command language interpreter rather than that of the invoked utility. The alternative
 9487 for such applications is the use of the *exec* family.

9488 C.2 Shell Command Language

9489 C.2.1 Shell Introduction

9490 The System V shell was selected as the starting point for the Shell and Utilities volume of
 9491 IEEE Std 1003.1-2001. The BSD C shell was excluded from consideration for the following
 9492 reasons:

- 9493 • Most historically portable shell scripts assume the Version 7 Bourne shell, from which the
 9494 System V shell is derived.
- 9495 • The majority of tutorial materials on shell programming assume the System V shell.

9496 The construct “#!” is reserved for implementations wishing to provide that extension. If it were
 9497 not reserved, the Shell and Utilities volume of IEEE Std 1003.1-2001 would disallow it by forcing
 9498 it to be a comment. As it stands, a strictly conforming application must not use “#!” as the first
 9499 two characters of the file.

9500 **C.2.2 Quoting**

9501 There is no additional rationale provided for this section.

9502 *C.2.2.1 Escape Character (Backslash)*

9503 There is no additional rationale provided for this section.

9504 *C.2.2.2 Single-Quotes*

9505 A backslash cannot be used to escape a single-quote in a single-quoted string. An embedded
 9506 quote can be created by writing, for example: "'a'\ 'b'", which yields "a'b". (See the Shell
 9507 and Utilities volume of IEEE Std 1003.1-2001, Section 2.6.5, Field Splitting for a better
 9508 understanding of how portions of words are either split into fields or remain concatenated.) A
 9509 single token can be made up of concatenated partial strings containing all three kinds of quoting
 9510 or escaping, thus permitting any combination of characters.

9511 *C.2.2.3 Double-Quotes*

9512 The escaped <newline> used for line continuation is removed entirely from the input and is not
 9513 replaced by any white space. Therefore, it cannot serve as a token separator.

9514 In double-quoting, if a backslash is immediately followed by a character that would be
 9515 interpreted as having a special meaning, the backslash is deleted and the subsequent character is
 9516 taken literally. If a backslash does not precede a character that would have a special meaning, it
 9517 is left in place unmodified and the character immediately following it is also left unmodified.
 9518 Thus, for example:

9519 "\\$" → \$

9520 "\a" → \a

9521 It would be desirable to include the statement "The characters from an enclosed "\${" to the
 9522 matching '}' shall not be affected by the double quotes", similar to the one for "\$()".
 9523 However, historical practice in the System V shell prevents this.

9524 The requirement that double-quotes be matched inside "\${...}" within double-quotes and the
 9525 rule for finding the matching '}' in the Shell and Utilities volume of IEEE Std 1003.1-2001,
 9526 Section 2.6.2, Parameter Expansion eliminate several subtle inconsistencies in expansion for
 9527 historical shells in rare cases; for example:

9528 "\${foo-bar}"

9529 yields **bar** when **foo** is not defined, and is an invalid substitution when **foo** is defined, in many
 9530 historical shells. The differences in processing the "\${...}" form have led to inconsistencies
 9531 between historical systems. A consequence of this rule is that single-quotes cannot be used to
 9532 quote the '}' within "\${...}"; for example:

9533 unset bar
 9534 foo="\${bar-'}'"

9535 is invalid because the "\${...}" substitution contains an unpaired unescaped single-quote. The
 9536 backslash can be used to escape the '}' in this example to achieve the desired result:

9537 unset bar
 9538 foo="\${bar-\\}'"

9539 The differences in processing the "\${...}" form have led to inconsistencies between the
 9540 historical System V shell, BSD, and KornShells, and the text in the Shell and Utilities volume of
 9541 IEEE Std 1003.1-2001 is an attempt to converge them without breaking too many applications.

9542 The only alternative to this compromise between shells would be to make the behavior
 9543 unspecified whenever the literal characters `'`, `{`, `}`, and `"` appear within `"${...}"`.
 9544 To write a portable script that uses these values, a user would have to assign variables; for
 9545 example:

```
9546     quote=\' dquote=\" lbrace='{ rbrace='}'
9547     ${foo-$quote$rbrace$quote}
```

9548 rather than:

```
9549     ${foo-"'"}
```

9550 Some implementations have allowed the end of the word to terminate the backquoted command
 9551 substitution, such as in:

```
9552     "`echo hello"
```

9553 This usage is undefined; the matching backquote is required by the Shell and Utilities volume of
 9554 IEEE Std 1003.1-2001. The other undefined usage can be illustrated by the example:

```
9555     sh -c '` echo "foo`'
```

9556 The description of the recursive actions involving command substitution can be illustrated with
 9557 an example. Upon recognizing the introduction of command substitution, the shell parses input
 9558 (in a new context), gathering the source for the command substitution until an unbalanced `'`
 9559 or `"` is located. For example, in the following:

```
9560     echo "$(date; echo "  

9561         one" )"
```

9562 the double-quote following the *echo* does not terminate the first double-quote; it is part of the
 9563 command substitution script. Similarly, in:

```
9564     echo "$(echo *)"
```

9565 the asterisk is not quoted since it is inside command substitution; however:

```
9566     echo "$(echo "*" )"
```

9567 is quoted (and represents the asterisk character itself).

9568 C.2.3 Token Recognition

9569 The `((` and `)` symbols are control operators in the KornShell, used for an alternative
 9570 syntax of an arithmetic expression command. A conforming application cannot use `((` as a
 9571 single token (with the exception of the `$(` form for shell arithmetic).

9572 On some implementations, the symbol `((` is a control operator; its use produces unspecified
 9573 results. Applications that wish to have nested subshells, such as:

```
9574     ((echo Hello);(echo World))
```

9575 must separate the `((` characters into two tokens by including white space between them.
 9576 Some systems may treat these as invalid arithmetic expressions instead of subshells.

9577 Certain combinations of characters are invalid in portable scripts, as shown in the grammar.
 9578 Implementations may use these combinations (such as `|&`) as valid control operators. Portable
 9579 scripts cannot rely on receiving errors in all cases where this volume of IEEE Std 1003.1-2001
 9580 indicates that a syntax is invalid.

9581 The (3) rule about combining characters to form operators is not meant to preclude systems from
 9582 extending the shell language when characters are combined in otherwise invalid ways.

9583 Conforming applications cannot use invalid combinations, and test suites should not penalize
 9584 systems that take advantage of this fact. For example, the unquoted combination "|&" is not
 9585 valid in a POSIX script, but has a specific KornShell meaning.

9586 The (10) rule about '#' as the current character is the first in the sequence in which a new token
 9587 is being assembled. The '#' starts a comment only when it is at the beginning of a token. This
 9588 rule is also written to indicate that the search for the end-of-comment does not consider escaped
 9589 <newline> specially, so that a comment cannot be continued to the next line.

9590 C.2.3.1 Alias Substitution

9591 The alias capability was added in the User Portability Utilities option because it is widely used in
 9592 historical implementations by interactive users.

9593 The definition of "alias name" precludes an alias name containing a slash character. Since the
 9594 text applies to the command words of simple commands, reserved words (in their proper
 9595 places) cannot be confused with aliases.

9596 The placement of alias substitution in token recognition makes it clear that it precedes all of the
 9597 word expansion steps.

9598 An example concerning trailing <blank>s and reserved words follows. If the user types:

```
9599     $ alias foo="/bin/ls "  
9600     $ alias while="/"
```

9601 The effect of executing:

```
9602     $ while true  
9603     > do  
9604     > echo "Hello, World"  
9605     > done
```

9606 is a never-ending sequence of "Hello, World" strings to the screen. However, if the user
 9607 types:

```
9608     $ foo while
```

9609 the result is an *ls* listing of /. Since the alias substitution for **foo** ends in a <space>, the next word
 9610 is checked for alias substitution. The next word, **while**, has also been aliased, so it is substituted
 9611 as well. Since it is not in the proper position as a command word, it is not recognized as a
 9612 reserved word.

9613 If the user types:

```
9614     $ foo; while
```

9615 **while** retains its normal reserved-word properties.

9616 C.2.4 Reserved Words

9617 All reserved words are recognized syntactically as such in the contexts described. However, note
 9618 that **in** is the only meaningful reserved word after a **case** or **for**; similarly, **in** is not meaningful as
 9619 the first word of a simple command.

9620 Reserved words are recognized only when they are delimited (that is, meet the definition of the
 9621 Base Definitions volume of IEEE Std 1003.1-2001, Section 3.435, Word), whereas operators are
 9622 themselves delimiters. For instance, '(' and ')' are control operators, so that no <space> is
 9623 needed in (*list*). However, '{' and '}' are reserved words in {*list*}, so that in this case the
 9624 leading <space> and semicolon are required.

9625 The list of unspecified reserved words is from the KornShell, so conforming applications cannot
 9626 use them in places a reserved word would be recognized. This list contained **time** in early
 9627 proposals, but it was removed when the *time* utility was selected for the Shell and Utilities
 9628 volume of IEEE Std 1003.1-2001.

9629 There was a strong argument for promoting braces to operators (instead of reserved words), so
 9630 they would be syntactically equivalent to subshell operators. Concerns about compatibility
 9631 outweighed the advantages of this approach. Nevertheless, conforming applications should
 9632 consider quoting '{' and '}' when they represent themselves.

9633 The restriction on ending a name with a colon is to allow future implementations that support
 9634 named labels for flow control; see the RATIONALE for the *break* built-in utility.

9635 It is possible that a future version of the Shell and Utilities volume of IEEE Std 1003.1-2001 may
 9636 require that '{' and '}' be treated individually as control operators, although the token "{}"
 9637 will probably be a special-case exemption from this because of the often-used *find*{ } construct.

9638 C.2.5 Parameters and Variables

9639 C.2.5.1 Positional Parameters

9640 There is no additional rationale provided for this section.

9641 C.2.5.2 Special Parameters

9642 Most historical implementations implement subshells by forking; thus, the special parameter
 9643 '\$' does not necessarily represent the process ID of the shell process executing the commands
 9644 since the subshell execution environment preserves the value of '\$'.

9645 If a subshell were to execute a background command, the value of "\$!" for the parent would
 9646 not change. For example:

```
9647 (
9648   date &
9649   echo $!
9650 )
9651 echo $!
```

9652 would echo two different values for "\$!".

9653 The "\$-" special parameter can be used to save and restore *set* options:

```
9654   Save=$(echo $- | sed 's/[ics]//g')
9655   ...
9656   set +aCefnuvx
9657   if [ -n "$Save" ]; then
9658       set -$Save
9659   fi
```

9660 The three options are removed using *sed* in the example because they may appear in the value of
 9661 "\$-" (from the *sh* command line), but are not valid options to *set*.

9662 The descriptions of parameters '*' and '@' assume the reader is familiar with the field splitting
 9663 discussion in the Shell and Utilities volume of IEEE Std 1003.1-2001, Section 2.6.5, Field Splitting
 9664 and understands that portions of the word remain concatenated unless there is some reason to
 9665 split them into separate fields.

9666 Some examples of the '*' and '@' properties, including the concatenation aspects:

```

9667      set "abc" "def ghi" "jkl"
9668      echo $*      => "abc" "def" "ghi" "jkl"
9669      echo "$*"    => "abc def ghi jkl"
9670      echo $@      => "abc" "def" "ghi" "jkl"

```

9671 but:

```

9672      echo "$@"      => "abc" "def ghi" "jkl"
9673      echo "xx$@yy"  => "xxabc" "def ghi" "jklyy"
9674      echo "$@$@"    => "abc" "def ghi" "jklabc" "def ghi" "jkl"

```

9675 In the preceding examples, the double-quote characters that appear after the "=>" do not appear
 9676 in the output and are used only to illustrate word boundaries.

9677 The following example illustrates the effect of setting *IFS* to a null string:

```

9678      $ IFS=' '
9679      $ set foo bar bam
9680      $ echo "$@"
9681      foo bar bam
9682      $ echo "$*"
9683      foobarbam
9684      $ unset IFS
9685      $ echo "$*"
9686      foo bar bam

```

9687 C.2.5.3 Shell Variables

9688 See the discussion of *IFS* in Section C.2.6.5 (on page 241) and the RATIONALE for the *sh* utility.

9689 The prohibition on *LC_CTYPE* changes affecting lexical processing protects the shell
 9690 implementor (and the shell programmer) from the ill effects of changing the definition of
 9691 <blank> or the set of alphabetic characters in the current environment. It would probably not be
 9692 feasible to write a compiled version of a shell script without this rule. The rule applies only to
 9693 the current invocation of the shell and its subshells—invoking a shell script or performing *exec sh*
 9694 would subject the new shell to the changes in *LC_CTYPE*.

9695 Other common environment variables used by historical shells are not specified by the Shell and
 9696 Utilities volume of IEEE Std 1003.1-2001, but they should be reserved for the historical uses.

9697 Tilde expansion for components of *PATH* in an assignment such as:

```

9698      PATH=~hlj/bin:~dwc/bin:$PATH

```

9699 is a feature of some historical shells and is allowed by the wording of the Shell and Utilities
 9700 volume of IEEE Std 1003.1-2001, Section 2.6.1, Tilde Expansion. Note that the tildes are expanded
 9701 during the assignment to *PATH*, not when *PATH* is accessed during command search.

9702 The following entries represent additional information about variables included in the Shell and
 9703 Utilities volume of IEEE Std 1003.1-2001, or rationale for common variables in use by shells that
 9704 have been excluded:

9705 — (Underscore.) While underscore is historical practice, its overloaded usage in
 9706 the KornShell is confusing, and it has been omitted from the Shell and Utilities
 9707 volume of IEEE Std 1003.1-2001.

9708 *ENV* This variable can be used to set aliases and other items local to the invocation
 9709 of a shell. The file referred to by *ENV* differs from *\$HOME/.profile* in that
 9710 *.profile* is typically executed at session start-up, whereas the *ENV* file is

9711		executed at the beginning of each shell invocation. The <i>ENV</i> value is
9712		interpreted in a manner similar to a dot script, in that the commands are
9713		executed in the current environment and the file needs to be readable, but not
9714		executable. However, unlike dot scripts, no <i>PATH</i> searching is performed.
9715		This is used as a guard against Trojan Horse security breaches.
9716	<i>ERRNO</i>	This variable was omitted from the Shell and Utilities volume of
9717		IEEE Std 1003.1-2001 because the values of error numbers are not defined in
9718		IEEE Std 1003.1-2001 in a portable manner.
9719	<i>FCEDIT</i>	Since this variable affects only the <i>fc</i> utility, it has been omitted from this more
9720		global place. The value of <i>FCEDIT</i> does not affect the command-line editing
9721		mode in the shell; see the description of <i>set -o vi</i> in the <i>set</i> built-in utility.
9722	<i>PS1</i>	This variable is used for interactive prompts. Historically, the “superuser”
9723		has had a prompt of '#'. Since privileges are not required to be monolithic, it
9724		is difficult to define which privileges should cause the alternate prompt.
9725		However, a sufficiently powerful user should be reminded of that power by
9726		having an alternate prompt.
9727	<i>PS3</i>	This variable is used by the KornShell for the <i>select</i> command. Since the POSIX
9728		shell does not include <i>select</i> , <i>PS3</i> was omitted.
9729	<i>PS4</i>	This variable is used for shell debugging. For example, the following script:
9730		<pre>PS4=' [\${LINENO}]+ '</pre>
9731		<pre>set -x</pre>
9732		<pre>echo Hello</pre>
9733		writes the following to standard error:
9734		<pre>[3]+ echo Hello</pre>
9735	<i>RANDOM</i>	This pseudo-random number generator was not seen as being useful to
9736		interactive users.
9737	<i>SECONDS</i>	Although this variable is sometimes used with <i>PS1</i> to allow the display of the
9738		current time in the prompt of the user, it is not one that would be manipulated
9739		frequently enough by an interactive user to include in the Shell and Utilities
9740		volume of IEEE Std 1003.1-2001.

9741 C.2.6 Word Expansions

9742 Step (2) refers to the “portions of fields generated by step (1)”. For example, if the word being
 9743 expanded were "\$x+\$y" and *IFS*=+, the word would be split only if "\$x" or "\$y" contained
 9744 '+'; the '+' in the original word was not generated by step (1).

9745 *IFS* is used for performing field splitting on the results of parameter and command substitution;
 9746 it is not used for splitting all fields. Previous versions of the shell used it for splitting all fields
 9747 during field splitting, but this has severe problems because the shell can no longer parse its own
 9748 script. There are also important security implications caused by this behavior. All useful
 9749 applications of *IFS* use it for parsing input of the *read* utility and for splitting the results of
 9750 parameter and command substitution.

9751 The rule concerning expansion to a single field requires that if **foo=abc** and **bar=def**, that:

9752 "\$foo" "\$bar"

9753 expands to the single field:

9754 abcdef

9755 The rule concerning empty fields can be illustrated by:

```
9756        $     unset foo
9757        $     set $foo bar ' ' xyz "$foo" abc
9758        $     for i
9759        >     do
9760        >         echo "-$i-"
9761        >     done
9762        -bar-
9763        --
9764        -xyz-
9765        --
9766        -abc-
```

9767 Step (1) indicates that parameter expansion, command substitution, and arithmetic expansion
9768 are all processed simultaneously as they are scanned. For example, the following is valid
9769 arithmetic:

```
9770        x=1
9771        echo $(( $(echo 3)+$x ))
```

9772 An early proposal stated that tilde expansion preceded the other steps, but this is not the case in
9773 known historical implementations; if it were, and if a referenced home directory contained a '\$'
9774 character, expansions would result within the directory name.

9775 *C.2.6.1 Tilde Expansion*

9776 Tilde expansion generally occurs only at the beginning of words, but an exception based on
9777 historical practice has been included:

```
9778        PATH=/posix/bin:~djk/bin
```

9779 This is eligible for tilde expansion because tilde follows a colon and none of the relevant
9780 characters is quoted. Consideration was given to prohibiting this behavior because any of the
9781 following are reasonable substitutes:

```
9782        PATH=$(printf %s ~karels/bin : ~bostic/bin)
9783        for Dir in ~maat/bin ~srb/bin ...
9784        do
9785        PATH=${PATH:+$PATH:}$Dir
9786        done
```

9787 In the first command, explicit colons are used for each directory. In all cases, the shell performs
9788 tilde expansion on each directory because all are separate words to the shell.

9789 Note that expressions in operands such as:

```
9790        make -k mumble LIBDIR=~chet/lib
```

9791 do not qualify as shell variable assignments, and tilde expansion is not performed (unless the
9792 command does so itself, which *make* does not).

9793 Because of the requirement that the word is not quoted, the following are not equivalent; only
9794 the last causes tilde expansion:

9795 \`~hlj/` \`~h\lj/` \`~"hlj"/` \`~hlj\` \`~hlj/`

9796 In an early proposal, tilde expansion occurred following any unquoted equals sign or colon, but
9797 this was removed because of its complexity and to avoid breaking commands such as:

9798 `rcp hostname:~marc/.profile .`

9799 A suggestion was made that the special sequence "`$~`" should be allowed to force tilde
9800 expansion anywhere. Since this is not historical practice, it has been left for future
9801 implementations to evaluate. (The description in the Shell and Utilities volume of
9802 IEEE Std 1003.1-2001, Section 2.2, Quoting requires that a dollar sign be quoted to represent
9803 itself, so the "`$~`" combination is already unspecified.)

9804 The results of giving tilde with an unknown login name are undefined because the KornShell
9805 "`~+`" and "`~-`" constructs make use of this condition, but in general it is an error to give an
9806 incorrect login name with tilde. The results of having *HOME* unset are unspecified because some
9807 historical shells treat this as an error.

9808 *C.2.6.2 Parameter Expansion*

9809 The rule for finding the closing '`'`' in "`${...}`" is the one used in the KornShell and is
9810 upwardly-compatible with the Bourne shell, which does not determine the closing '`'`' until the
9811 word is expanded. The advantage of this is that incomplete expansions, such as:

9812 `${foo`

9813 can be determined during tokenization, rather than during expansion.

9814 The string length and substring capabilities were included because of the demonstrated need for
9815 them, based on their usage in other shells, such as C shell and KornShell.

9816 Historical versions of the KornShell have not performed tilde expansion on the word part of
9817 parameter expansion; however, it is more consistent to do so.

9818 *C.2.6.3 Command Substitution*

9819 The "`$()`" form of command substitution solves a problem of inconsistent behavior when using
9820 backquotes. For example:

9821	Command	Output
9822	<code>echo '\\$x'</code>	<code>\\$x</code>
9823	<code>echo `echo '\\$x'`</code>	<code>\$x</code>
9824	<code>echo \$(echo '\\$x')</code>	<code>\\$x</code>

9825 Additionally, the backquoted syntax has historical restrictions on the contents of the embedded
9826 command. While the newer "`$()`" form can process any kind of valid embedded script, the
9827 backquoted form cannot handle some valid scripts that include backquotes. For example, these
9828 otherwise valid embedded scripts do not work in the left column, but do work on the right:


```

9829         echo `          echo $(
9830         cat <<\eof        cat <<\eof
9831         a here-doc with `  a here-doc with )
9832         eof              eof
9833         `                )

9834         echo `          echo $(
9835         echo abc # a comment with `  echo abc # a comment with )
9836         `                )

9837         echo `          echo $(
9838         echo ` ` `        echo ` ` `
9839         `                )

```

9840 Because of these inconsistent behaviors, the backquoted variety of command substitution is not
 9841 recommended for new applications that nest command substitutions or attempt to embed
 9842 complex scripts.

9843 The KornShell feature:

9844 If *command* is of the form `<word`, *word* is expanded to generate a pathname, and the value of
 9845 the command substitution is the contents of this file with any trailing `<newline>`s deleted.

9846 was omitted from the Shell and Utilities volume of IEEE Std 1003.1-2001 because `$(cat word)` is
 9847 an appropriate substitute. However, to prevent breaking numerous scripts relying on this
 9848 feature, it is unspecified to have a script within `"$()"` that has only redirections.

9849 The requirement to separate `"$("` and `' ('` when a single subshell is command-substituted is to
 9850 avoid any ambiguities with arithmetic expansion.

9851 C.2.6.4 Arithmetic Expansion

9852 The `"(())"` form of KornShell arithmetic in early proposals was omitted. The standard
 9853 developers concluded that there was a strong desire for some kind of arithmetic evaluator to
 9854 replace *expr*, and that relating it to `'$'` makes it work well with the standard shell language, and
 9855 it provides access to arithmetic evaluation in places where accessing a utility would be
 9856 inconvenient.

9857 The syntax and semantics for arithmetic were changed for the ISO/IEC 9945-2:1993 standard.
 9858 The language is essentially a pure arithmetic evaluator of constants and operators (excluding
 9859 assignment) and represents a simple subset of the previous arithmetic language (which was
 9860 derived from the KornShell `"(())"` construct). The syntax was changed from that of a
 9861 command denoted by `((expression))` to an expansion denoted by `$(expression)`. The new form is
 9862 a dollar expansion `'$'` that evaluates the expression and substitutes the resulting value.
 9863 Objections to the previous style of arithmetic included that it was too complicated, did not fit in
 9864 well with the use of variables in the shell, and its syntax conflicted with subshells. The
 9865 justification for the new syntax is that the shell is traditionally a macro language, and if a new
 9866 feature is to be added, it should be accomplished by extending the capabilities presented by the
 9867 current model of the shell, rather than by inventing a new one outside the model; adding a new
 9868 dollar expansion was perceived to be the most intuitive and least destructive way to add such a
 9869 new capability.

9870 In early proposals, a form `$(expression)` was used. It was functionally equivalent to the `"$(())"`
 9871 of the current text, but objections were lodged that the 1988 KornShell had already implemented
 9872 `"$(())"` and there was no compelling reason to invent yet another syntax. Furthermore, the
 9873 `"$[]"` syntax had a minor incompatibility involving the patterns in **case** statements.

9874 The portion of the ISO C standard arithmetic operations selected corresponds to the operations
9875 historically supported in the KornShell.

9876 It was concluded that the *test* command (D) was sufficient for the majority of relational arithmetic
9877 tests, and that tests involving complicated relational expressions within the shell are rare, yet
9878 could still be accommodated by testing the value of "\$(())" itself. For example:

```
9879     # a complicated relational expression
9880     while [ $(( (($x + $y)/($a * $b)) < ($foo*$bar) )) -ne 0 ]
```

9881 or better yet, the rare script that has many complex relational expressions could define a
9882 function like this:

```
9883     val() {
9884         return $(!$1)
9885     }
```

9886 and complicated tests would be less intimidating:

```
9887     while val $(( (($x + $y)/($a * $b)) < ($foo*$bar) ))
9888     do
9889         # some calculations
9890     done
```

9891 A suggestion that was not adopted was to modify *true* and *false* to take an optional argument,
9892 and *true* would exit true only if the argument was non-zero, and *false* would exit false only if the
9893 argument was non-zero:

```
9894     while true $((($x > 5 && $y <= 25))
```

9895 There is a minor portability concern with the new syntax. The example "\$((2+2))" could have
9896 been intended to mean a command substitution of a utility named "2+2" in a subshell. The
9897 standard developers considered this to be obscure and isolated to some KornShell scripts
9898 (because "\$()" command substitution existed previously only in the KornShell). The text on
9899 command substitution requires that the "\$(" and '((' be separate tokens if this usage is needed.

9900 An example such as:

```
9901     echo $((echo hi);(echo there))
```

9902 should not be misinterpreted by the shell as arithmetic because attempts to balance the
9903 parentheses pairs would indicate that they are subshells. However, as indicated by the Base
9904 Definitions volume of IEEE Std 1003.1-2001, Section 3.112, Control Operator, a conforming
9905 application must separate two adjacent parentheses with white space to indicate nested
9906 subshells.

9907 Although the ISO/IEC 9899:1999 standard now requires support for **long long** and allows
9908 extended integer types with higher ranks, IEEE Std 1003.1-2001 only requires arithmetic
9909 expansions to support **signed long** integer arithmetic. Implementations are encouraged to
9910 support signed integer values at least as large as the size of the largest file allowed on the
9911 implementation.

9912 Implementations are also allowed to perform floating-point evaluations as long as an
9913 application won't see different results for expressions that would not overflow **signed long**
9914 integer expression evaluation. (This includes appropriate truncation of results to integer values.)

9915 Changes made in response to IEEE PASC Interpretation 1003.2 #208 removed the requirement
9916 that the integer constant suffixes **l** and **L** had to be recognized. The ISO POSIX-2:1993 standard
9917 did not require the **u**, **uL**, **uL**, **U**, **UL**, **UL**, **lu**, **lU**, **Lu**, and **LU** suffixes since only signed integer
9918 arithmetic was required. Since all arithmetic expressions were treated as handling **signed long**

9919 integer types anyway, the `l` and `L` suffixes were redundant. No known scripts used them and
 9920 some historic shells did not support them. When the ISO/IEC 9899:1999 standard was used as
 9921 the basis for the description of arithmetic processing, the `ll` and `LL` suffixes and combinations
 9922 were also not required. Implementations are still free to accept any or all of these suffixes, but
 9923 are not required to do so.

9924 There was also some confusion as to whether the shell was required to recognize character
 9925 constants. Syntactically, character constants were required to be recognized, but the
 9926 requirements for the handling of backslash (`'\'`) and quote (`' ''`) characters (needed to specify
 9927 character constants) within an arithmetic expansion were ambiguous. Furthermore, no known
 9928 shells supported them. Changes made in response to IEEE PASC Interpretation 1003.2 #208
 9929 removed the requirement to support them (if they were indeed required before).
 9930 IEEE Std 1003.1-2001 clearly does not require support for character constants.

9931 *C.2.6.5 Field Splitting*

9932 The operation of field splitting using *IFS*, as described in early proposals, was based on the way
 9933 the KornShell splits words, but it is incompatible with other common versions of the shell.
 9934 However, each has merit, and so a decision was made to allow both. If the *IFS* variable is unset
 9935 or is `<space><tab><newline>`, the operation is equivalent to the way the System V shell splits
 9936 words. Using characters outside the `<space><tab><newline>` set yields the KornShell behavior,
 9937 where each of the non-`<space><tab><newline>`s is significant. This behavior, which affords the
 9938 most flexibility, was taken from the way the original *awk* handled field splitting.

9939 Rule (3) can be summarized as a pseudo-ERE:

```
9940 (s*ns*|s+)
```

9941 where *s* is an *IFS* white space character and *n* is a character in the *IFS* that is not white space.
 9942 Any string matching that ERE delimits a field, except that the *s+* form does not delimit fields at
 9943 the beginning or the end of a line. For example, if *IFS* is `<space>/<comma>/<tab>`, the string:

```
9944 <space><space>red<space><space>, <space>white<space>blue
```

9945 yields the three colors as the delimited fields.

9946 *C.2.6.6 Pathname Expansion*

9947 There is no additional rationale provided for this section.

9948 *C.2.6.7 Quote Removal*

9949 There is no additional rationale provided for this section.

9950 **C.2.7 Redirection**

9951 In the System Interfaces volume of IEEE Std 1003.1-2001, file descriptors are integers in the range
 9952 `0-({OPEN_MAX}-1)`. The file descriptors discussed in the Shell and Utilities volume of
 9953 IEEE Std 1003.1-2001, Section 2.7, Redirection are that same set of small integers.

9954 Having multi-digit file descriptor numbers for I/O redirection can cause some obscure
 9955 compatibility problems. Specifically, scripts that depend on an example command:

```
9956 echo 22>/dev/null
```

9957 echoing `"2"` to standard error or `"22"` to standard output are no longer portable. However, the
 9958 file descriptor number must still be delimited from the preceding text. For example:

9959 cat file2>foo

9960 writes the contents of **file2**, not the contents of **file**.

9961 The ">|" format of output redirection was adopted from the KornShell. Along with the
9962 *noclobber* option, *set -C*, it provides a safety feature to prevent inadvertent overwriting of
9963 existing files. (See the RATIONALE for the *pathchk* utility for why this step was taken.) The
9964 restriction on regular files is historical practice.

9965 The System V shell and the KornShell have differed historically on pathname expansion of *word*;
9966 the former never performed it, the latter only when the result was a single field (file). As a
9967 compromise, it was decided that the KornShell functionality was useful, but only as a shorthand
9968 device for interactive users. No reasonable shell script would be written with a command such
9969 as:

9970 cat foo > a*

9971 Thus, shell scripts are prohibited from doing it, while interactive users can select the shell with
9972 which they are most comfortable.

9973 The construct "2>&1" is often used to redirect standard error to the same file as standard
9974 output. Since the redirections take place beginning to end, the order of redirections is significant.
9975 For example:

9976 ls > foo 2>&1

9977 directs both standard output and standard error to file **foo**. However:

9978 ls 2>&1 > foo

9979 only directs standard output to file **foo** because standard error was duplicated as standard
9980 output before standard output was directed to file **foo**.

9981 The "<>" operator could be useful in writing an application that worked with several terminals,
9982 and occasionally wanted to start up a shell. That shell would in turn be unable to run
9983 applications that run from an ordinary controlling terminal unless it could make use of "<>"
9984 redirection. The specific example is a historical version of the pager *more*, which reads from
9985 standard error to get its commands, so standard input and standard output are both available
9986 for their usual usage. There is no way of saying the following in the shell without "<>":

9987 cat food | more - >/dev/tty03 2<>/dev/tty03

9988 Another example of "<>" is one that opens **/dev/tty** on file descriptor 3 for reading and writing:

9989 exec 3<> /dev/tty

9990 An example of creating a lock file for a critical code region:

```
9991           set -C
9992           until    2> /dev/null > lockfile
9993           do       sleep 30
9994           done
9995           set +C
9996           perform critical function
9997           rm lockfile
```

9998 Since **/dev/null** is not a regular file, no error is generated by redirecting to it in *noclobber* mode.

9999 Tilde expansion is not performed on a here-document because the data is treated as if it were
10000 enclosed in double quotes.

10001 *C.2.7.1 Redirecting Input*

10002 There is no additional rationale provided for this section.

10003 *C.2.7.2 Redirecting Output*

10004 There is no additional rationale provided for this section.

10005 *C.2.7.3 Appending Redirected Output*

10006 Note that when a file is opened (even with the `O_APPEND` flag set), the initial file offset for that
10007 file is set to the beginning of the file. Some historic shells set the file offset to the current end-of-
10008 file when **append** mode shell redirection was used, but this is not allowed by
10009 IEEE Std 1003.1-2001.

10010 *C.2.7.4 Here-Document*

10011 There is no additional rationale provided for this section.

10012 *C.2.7.5 Duplicating an Input File Descriptor*

10013 There is no additional rationale provided for this section.

10014 *C.2.7.6 Duplicating an Output File Descriptor*

10015 There is no additional rationale provided for this section.

10016 *C.2.7.7 Open File Descriptors for Reading and Writing*

10017 There is no additional rationale provided for this section.

10018 **C.2.8 Exit Status and Errors**10019 *C.2.8.1 Consequences of Shell Errors*

10020 There is no additional rationale provided for this section.

10021 *C.2.8.2 Exit Status for Commands*

10022 There is a historical difference in *sh* and *ksh* non-interactive error behavior. When a command
10023 named in a script is not found, some implementations of *sh* exit immediately, but *ksh* continues
10024 with the next command. Thus, the Shell and Utilities volume of IEEE Std 1003.1-2001 says that
10025 the shell “may” exit in this case. This puts a small burden on the programmer, who has to test
10026 for successful completion following a command if it is important that the next command not be
10027 executed if the previous command was not found. If it is important for the command to have
10028 been found, it was probably also important for it to complete successfully. The test for successful
10029 completion would not need to change.

10030 Historically, shells have returned an exit status of $128+n$, where n represents the signal number.
10031 Since signal numbers are not standardized, there is no portable way to determine which signal
10032 caused the termination. Also, it is possible for a command to exit with a status in the same range
10033 of numbers that the shell would use to report that the command was terminated by a signal.
10034 Implementations are encouraged to choose exit values greater than 256 to indicate programs
10035 that terminate by a signal so that the exit status cannot be confused with an exit status generated
10036 by a normal termination.

10037 Historical shells make the distinction between “utility not found” and “utility found but cannot
10038 execute” in their error messages. By specifying two seldomly used exit status values for these

10039 cases, 127 and 126 respectively, this gives an application the opportunity to make use of this
 10040 distinction without having to parse an error message that would probably change from locale to
 10041 locale. The *command*, *env*, *nohup*, and *xargs* utilities in the Shell and Utilities volume of
 10042 IEEE Std 1003.1-2001 have also been specified to use this convention.

10043 When a command fails during word expansion or redirection, most historical implementations
 10044 exit with a status of 1. However, there was some sentiment that this value should probably be
 10045 much higher so that an application could distinguish this case from the more normal exit status
 10046 values. Thus, the language “greater than zero” was selected to allow either method to be
 10047 implemented.

10048 C.2.9 Shell Commands

10049 A description of an “empty command” was removed from an early proposal because it is only
 10050 relevant in the cases of *sh -c ""*, *system("")*, or an empty shell-script file (such as the
 10051 implementation of *true* on some historical systems). Since it is no longer mentioned in the Shell
 10052 and Utilities volume of IEEE Std 1003.1-2001, it falls into the silently unspecified category of
 10053 behavior where implementations can continue to operate as they have historically, but
 10054 conforming applications do not construct empty commands. (However, note that *sh* does
 10055 explicitly state an exit status for an empty string or file.) In an interactive session or a script with
 10056 other commands, extra <newline>s or semicolons, such as:

```
10057     $ false
10058     $
10059     $ echo $?
10060     1
```

10061 would not qualify as the empty command described here because they would be consumed by
 10062 other parts of the grammar.

10063 C.2.9.1 Simple Commands

10064 The enumerated list is used only when the command is actually going to be executed. For
 10065 example, in:

```
10066     true || $foo *
```

10067 no expansions are performed.

10068 The following example illustrates both how a variable assignment without a command name
 10069 affects the current execution environment, and how an assignment with a command name only
 10070 affects the execution environment of the command:

```
10071     $ x=red
10072     $ echo $x
10073     red
10074     $ export x
10075     $ sh -c 'echo $x'
10076     red
10077     $ x=blue sh -c 'echo $x'
10078     blue
10079     $ echo $x
10080     red
```

10081 This next example illustrates that redirections without a command name are still performed:

```

10082     $ ls foo
10083     ls: foo: no such file or directory
10084     $ > foo
10085     $ ls foo
10086     foo

```

10087 A command without a command name, but one that includes a command substitution, has an
 10088 exit status of the last command substitution that the shell performed. For example:

```

10089     if      x=$(command)
10090     then    ...
10091     fi

```

10092 An example of redirections without a command name being performed in a subshell shows that
 10093 the here-document does not disrupt the standard input of the **while** loop:

```

10094     IFS=:
10095     while  read a b
10096     do    echo $a
10097           <<-eof
10098           Hello
10099           eof
10100     done </etc/passwd

```

10101 Following are examples of commands without command names in AND-OR lists:

```

10102     > foo || {
10103         echo "error: foo cannot be created" >&2
10104         exit 1
10105     }
10106     # set saved if /vmunix.save exists
10107     test -f /vmunix.save && saved=1

```

10108 Command substitution and redirections without command names both occur in subshells, but
 10109 they are not necessarily the same ones. For example, in:

```

10110     exec 3> file
10111     var=$(echo foo >&3) 3>&1

```

10112 it is unspecified whether **foo** is echoed to the file or to standard output.

10113 Command Search and Execution

10114 This description requires that the shell can execute shell scripts directly, even if the underlying
 10115 system does not support the common "#!" interpreter convention. That is, if file **foo** contains
 10116 shell commands and is executable, the following executes **foo**:

```

10117     ./foo

```

10118 The command search shown here does not match all historical implementations. A more typical
 10119 sequence has been:

- 10120 • Any built-in (special or regular)
- 10121 • Functions
- 10122 • Path search for executable files

10123 But there are problems with this sequence. Since the programmer has no idea in advance which
 10124 utilities might have been built into the shell, a function cannot be used to override portably a

10125 utility of the same name. (For example, a function named *cd* cannot be written for many
10126 historical systems.) Furthermore, the *PATH* variable is partially ineffective in this case, and only
10127 a pathname with a slash can be used to ensure a specific executable file is invoked.

10128 After the *execve()* failure described, the shell normally executes the file as a shell script. Some
10129 implementations, however, attempt to detect whether the file is actually a script and not an
10130 executable from some other architecture. The method used by the KornShell is allowed by the
10131 text that indicates non-text files may be bypassed.

10132 The sequence selected for the Shell and Utilities volume of IEEE Std 1003.1-2001 acknowledges
10133 that special built-ins cannot be overridden, but gives the programmer full control over which
10134 versions of other utilities are executed. It provides a means of suppressing function lookup (via
10135 the *command* utility) for the user's own functions and ensures that any regular built-ins or
10136 functions provided by the implementation are under the control of the path search. The
10137 mechanisms for associating built-ins or functions with executable files in the path are not
10138 specified by the Shell and Utilities volume of IEEE Std 1003.1-2001, but the wording requires that
10139 if either is implemented, the application is not able to distinguish a function or built-in from an
10140 executable (other than in terms of performance, presumably). The implementation ensures that
10141 all effects specified by the Shell and Utilities volume of IEEE Std 1003.1-2001 resulting from the
10142 invocation of the regular built-in or function (interaction with the environment, variables, traps,
10143 and so on) are identical to those resulting from the invocation of an executable file.

10144 **Examples**

10145 Consider three versions of the *ls* utility:

- 10146 1. The application includes a shell function named *ls*.
- 10147 2. The user writes a utility named *ls* and puts it in **/fred/bin**.
- 10148 3. The example implementation provides *ls* as a regular shell built-in that is invoked (either
10149 by the shell or directly by *exec*) when the path search reaches the directory **/posix/bin**.

10150 If *PATH*=**/posix/bin**, various invocations yield different versions of *ls*:

10151

10152

10153

10154

10155

10156

10157

Invocation	Version of <i>ls</i>
<i>ls</i> (from within application script)	(1) function
<i>command ls</i> (from within application script)	(3) built-in
<i>ls</i> (from within makefile called by application)	(3) built-in
<i>system("ls")</i>	(3) built-in
<i>PATH="/fred/bin:\$PATH" ls</i>	(2) user's version

10158 *C.2.9.2 Pipelines*

10159 Because pipeline assignment of standard input or standard output or both takes place before
10160 redirection, it can be modified by redirection. For example:

```
10161 $ command1 2>&1 | command2
```

10162 sends both the standard output and standard error of *command1* to the standard input of
10163 *command2*.

10164 The reserved word **!** allows more flexible testing using AND and OR lists.

10165 It was suggested that it would be better to return a non-zero value if any command in the
10166 pipeline terminates with non-zero status (perhaps the bitwise-inclusive OR of all return values).
10167 However, the choice of the last-specified command semantics are historical practice and would

10168 cause applications to break if changed. An example of historical behavior:

```
10169     $ sleep 5 | (exit 4)
10170     $ echo $?
10171     4
10172     $ (exit 4) | sleep 5
10173     $ echo $?
10174     0
```

10175 C.2.9.3 Lists

10176 The equal precedence of "&&" and "||" is historical practice. The standard developers
 10177 evaluated the model used more frequently in high-level programming languages, such as C, to
 10178 allow the shell logical operators to be used for complex expressions in an unambiguous way, but
 10179 they could not allow historical scripts to break in the subtle way unequal precedence might
 10180 cause. Some arguments were posed concerning the "{" or "(" groupings that are required
 10181 historically. There are some disadvantages to these groupings:

- 10182 • The "(" can be expensive, as they spawn other processes on some implementations. This
 10183 performance concern is primarily an implementation issue.
- 10184 • The "{" braces are not operators (they are reserved words) and require a trailing space
 10185 after each '{', and a semicolon before each '}'. Most programmers (and certainly
 10186 interactive users) have avoided braces as grouping constructs because of the problematic
 10187 syntax required. Braces were not changed to operators because that would generate
 10188 compatibility issues even greater than the precedence question; braces appear outside the
 10189 context of a keyword in many shell scripts.

10190 IEEE PASC Interpretation 1003.2 #204 is applied, clarifying that the operators "&&" and "||"
 10191 are evaluated with left associativity.

10192 Asynchronous Lists

10193 The grammar treats a construct such as:

```
10194     foo & bar & bam &
```

10195 as one "asynchronous list", but since the status of each element is tracked by the shell, the term
 10196 "element of an asynchronous list" was introduced to identify just one of the **foo**, **bar**, or **bam**
 10197 portions of the overall list.

10198 Unless the implementation has an internal limit, such as {CHILD_MAX}, on the retained process
 10199 IDs, it would require unbounded memory for the following example:

```
10200     while true
10201     do         foo & echo $!
10202     done
```

10203 The treatment of the signals SIGINT and SIGQUIT with asynchronous lists is described in the
 10204 Shell and Utilities volume of IEEE Std 1003.1-2001, Section 2.11, Signals and Error Handling.

10205 Since the connection of the input to the equivalent of /dev/null is considered to occur before
 10206 redirections, the following script would produce no output:

```
10207     exec < /etc/passwd
10208     cat <&0 &
10209     wait
```

- 10210 **Sequential Lists**
- 10211 There is no additional rationale provided for this section.
- 10212 **AND Lists**
- 10213 There is no additional rationale provided for this section.
- 10214 **OR Lists**
- 10215 There is no additional rationale provided for this section.
- 10216 *C.2.9.4 Compound Commands*
- 10217 **Grouping Commands**
- 10218 The semicolon shown in *{compound-list;}* is an example of a control operator delimiting the } reserved word. Other delimiters are possible, as shown in the Shell and Utilities volume of IEEE Std 1003.1-2001, Section 2.10, Shell Grammar; <newline> is frequently used.
- 10219
- 10220
- 10221 A proposal was made to use the <do-done> construct in all cases where command grouping in the current process environment is performed, identifying it as a construct for the grouping commands, as well as for shell functions. This was not included because the shell already has a grouping construct for this purpose ("{}"), and changing it would have been counter-productive.
- 10222
- 10223
- 10224
- 10225
- 10226 **For Loop**
- 10227 The format is shown with generous usage of <newline>s. See the grammar in the Shell and Utilities volume of IEEE Std 1003.1-2001, Section 2.10, Shell Grammar for a precise description of where <newline>s and semicolons can be interchanged.
- 10228
- 10229
- 10230 Some historical implementations support ' { ' and ' } ' as substitutes for **do** and **done**. The standard developers chose to omit them, even as an obsolescent feature. (Note that these substitutes were only for the **for** command; the **while** and **until** commands could not use them historically because they are followed by compound-lists that may contain "{ . . . }" grouping commands themselves.)
- 10231
- 10232
- 10233
- 10234
- 10235 The reserved word pair **do** ... **done** was selected rather than **do** ... **od** (which would have matched the spirit of **if** ... **fi** and **case** ... **esac**) because *od* is already the name of a standard utility.
- 10236
- 10237
- 10238 PASC Interpretation 1003.2 #169 has been applied changing the grammar.
- 10239 **Case Conditional Construct**
- 10240 An optional left parenthesis before *pattern* was added to allow numerous historical KornShell scripts to conform. At one time, using the leading parenthesis was required if the **case** statement was to be embedded within a "\$ ()" command substitution; this is no longer the case with the POSIX shell. Nevertheless, many historical scripts use the left parenthesis, if only because it makes matching-parenthesis searching easier in *vi* and other editors. This is a relatively simple implementation change that is upwards-compatible for all scripts.
- 10241
- 10242
- 10243
- 10244
- 10245
- 10246 Consideration was given to requiring *break* inside the *compound-list* to prevent falling through to the next pattern action list. This was rejected as being nonexistent practice. An interesting undocumented feature of the KornShell is that using "&" instead of ";" as a terminator causes the exact opposite behavior—the flow of control continues with the next *compound-list*.
- 10247
- 10248
- 10249

10250 The pattern ' * ', given as the last pattern in a **case** construct, is equivalent to the default case in
10251 a C-language **switch** statement.

10252 The grammar shows that reserved words can be used as patterns, even if one is the first word on
10253 a line. Obviously, the reserved word **esac** cannot be used in this manner.

10254 **If Conditional Construct**

10255 The precise format for the command syntax is described in the Shell and Utilities volume of
10256 IEEE Std 1003.1-2001, Section 2.10, Shell Grammar.

10257 **While Loop**

10258 The precise format for the command syntax is described in the Shell and Utilities volume of
10259 IEEE Std 1003.1-2001, Section 2.10, Shell Grammar.

10260 **Until Loop**

10261 The precise format for the command syntax is described in the Shell and Utilities volume of
10262 IEEE Std 1003.1-2001, Section 2.10, Shell Grammar.

10263 *C.2.9.5 Function Definition Command*

10264 The description of functions in an early proposal was based on the notion that functions should
10265 behave like miniature shell scripts; that is, except for sharing variables, most elements of an
10266 execution environment should behave as if they were a new execution environment, and
10267 changes to these should be local to the function. For example, traps and options should be reset
10268 on entry to the function, and any changes to them do not affect the traps or options of the caller.
10269 There were numerous objections to this basic idea, and the opponents asserted that functions
10270 were intended to be a convenient mechanism for grouping common commands that were to be
10271 executed in the current execution environment, similar to the execution of the *dot* special
10272 built-in.

10273 It was also pointed out that the functions described in that early proposal did not provide a local
10274 scope for everything a new shell script would, such as the current working directory, or *umask*,
10275 but instead provided a local scope for only a few select properties. The basic argument was that
10276 if a local scope is needed for the execution environment, the mechanism already existed: the
10277 application can put the commands in a new shell script and call that script. All historical shells
10278 that implemented functions, other than the KornShell, have implemented functions that operate
10279 in the current execution environment. Because of this, traps and options have a global scope
10280 within a shell script. Local variables within a function were considered and included in another
10281 early proposal (controlled by the special built-in *local*), but were removed because they do not fit
10282 the simple model developed for functions and because there was some opposition to adding yet
10283 another new special built-in that was not part of historical practice. Implementations should
10284 reserve the identifier *local* (as well as *typeset*, as used in the KornShell) in case this local variable
10285 mechanism is adopted in a future version of IEEE Std 1003.1-2001.

10286 A separate issue from the execution environment of a function is the availability of that function
10287 to child shells. A few objectors maintained that just as a variable can be shared with child shells
10288 by exporting it, so should a function. In early proposals, the *export* command therefore had a *-f*
10289 flag for exporting functions. Functions that were exported were to be put into the environment
10290 as *name()=value* pairs, and upon invocation, the shell would scan the environment for these and
10291 automatically define these functions. This facility was strongly opposed and was omitted. Some
10292 of the arguments against exportable functions were as follows:

- 10293 • There was little historical practice. The Ninth Edition shell provided them, but there was
10294 controversy over how well it worked.
- 10295 • There are numerous security problems associated with functions appearing in the
10296 environment of a user and overriding standard utilities or the utilities owned by the
10297 application.
- 10298 • There was controversy over requiring *make* to import functions, where it has historically used
10299 an *exec* function for many of its command line executions.
- 10300 • Functions can be big and the environment is of a limited size. (The counter-argument was
10301 that functions are no different from variables in terms of size: there can be big ones, and there
10302 can be small ones—and just as one does not export huge variables, one does not export huge
10303 functions. However, this might not apply to the average shell-function writer, who typically
10304 writes much larger functions than variables.)

10305 As far as can be determined, the functions in the Shell and Utilities volume of
10306 IEEE Std 1003.1-2001 match those in System V. Earlier versions of the KornShell had two
10307 methods of defining functions:

```
10308       function fname { compound-list }
```

10309 and:

```
10310       fname() { compound-list }
```

10311 The latter used the same definition as the Shell and Utilities volume of IEEE Std 1003.1-2001, but
10312 differed in semantics, as described previously. The current edition of the KornShell aligns the
10313 latter syntax with the Shell and Utilities volume of IEEE Std 1003.1-2001 and keeps the former as
10314 is.

10315 The name space for functions is limited to that of a *name* because of historical practice.
10316 Complications in defining the syntactic rules for the function definition command and in dealing
10317 with known extensions such as the "@()" usage in the KornShell prevented the name space
10318 from being widened to a *word*. Using functions to support synonyms such as the "!!" and '%'
10319 usage in the C shell is thus disallowed to conforming applications, but acceptable as an
10320 extension. For interactive users, the aliasing facilities in the Shell and Utilities volume of
10321 IEEE Std 1003.1-2001 should be adequate for this purpose. It is recognized that the name space
10322 for utilities in the file system is wider than that currently supported for functions, if the portable
10323 filename character set guidelines are ignored, but it did not seem useful to mandate extensions
10324 in systems for so little benefit to conforming applications.

10325 The "()" in the function definition command consists of two operators. Therefore, intermixing
10326 <blank>s with the *fname*, '(', and ')' is allowed, but unnecessary.

10327 An example of how a function definition can be used wherever a simple command is allowed:

```
10328       # If variable i is equal to "yes",
10329       # define function foo to be ls -l
10330       #
10331       [ "$i" = yes ] && foo() {
10332        ls -l
10333       }
```

10334 **C.2.10 Shell Grammar**

10335 There are several subtle aspects of this grammar where conventional usage implies rules about
10336 the grammar that in fact are not true.

10337 For *compound_list*, only the forms that end in a *separator* allow a reserved word to be recognized,
10338 so usually only a *separator* can be used where a compound list precedes a reserved word (such as
10339 **Then, Else, Do,** and **Rbrace**). Explicitly requiring a separator would disallow such valid (if rare)
10340 statements as:

```
10341     if (false) then (echo x) else (echo y) fi
```

10342 See the Note under special grammar rule (1).

10343 Concerning the third sentence of rule (1) (“Also, if the parser ...”):

10344 • This sentence applies rather narrowly: when a compound list is terminated by some clear
10345 delimiter (such as the closing **fi** of an inner **if_clause**) then it would apply; where the
10346 compound list might continue (as in after a ‘;’), rule (7a) (and consequently the first
10347 sentence of rule (1)) would apply. In many instances the two conditions are identical, but this
10348 part of rule (1) does not give license to treating a **WORD** as a reserved word unless it is in a
10349 place where a reserved word has to appear.

10350 • The statement is equivalent to requiring that when the LR(1) lookahead set contains exactly
10351 one reserved word, it must be recognized if it is present. (Here “LR(1)” refers to the
10352 theoretical concepts, not to any real parser generator.)

10353 For example, in the construct below, and when the parser is at the point marked with ‘^’,
10354 the only next legal token is **then** (this follows directly from the grammar rules):

```
10355     if if...fi then ... fi
10356         ^
```

10357 At that point, the **then** must be recognized as a reserved word.

10358 (Depending on the parser generator actually used, “extra” reserved words may be in some
10359 lookahead sets. It does not really matter if they are recognized, or even if any possible
10360 reserved word is recognized in that state, because if it is recognized and is not in the
10361 (theoretical) LR(1) lookahead set, an error is ultimately detected. In the example above, if
10362 some other reserved word (for example, **while**) is also recognized, an error occurs later.

10363 This is approximately equivalent to saying that reserved words are recognized after other
10364 reserved words (because it is after a reserved word that this condition occurs), but avoids the
10365 “except for ...” list that would be required for **case**, **for**, and so on. (Reserved words are of
10366 course recognized anywhere a *simple_command* can appear, as well. Other rules take care of
10367 the special cases of non-recognition, such as rule (4) for **case** statements.)

10368 Note that the body of here-documents are handled by token recognition (see the Shell and
10369 Utilities volume of IEEE Std 1003.1-2001, Section 2.3, Token Recognition) and do not appear in
10370 the grammar directly. (However, the here-document I/O redirection operator is handled as part
10371 of the grammar.)

10372 The start symbol of the grammar (**complete_command**) represents either input from the
10373 command line or a shell script. It is repeatedly applied by the interpreter to its input and
10374 represents a single “chunk” of that input as seen by the interpreter.

10375 *C.2.10.1 Shell Grammar Lexical Conventions*

10376 There is no additional rationale provided for this section.

10377 *C.2.10.2 Shell Grammar Rules*

10378 There is no additional rationale provided for this section.

10379 **C.2.11 Signals and Error Handling**

10380 There is no additional rationale provided for this section.

10381 **C.2.12 Shell Execution Environment**10382 Some implementations have implemented the last stage of a pipeline in the current environment
10383 so that commands such as:10384 `command | read foo`10385 set variable **foo** in the current environment. This extension is allowed, but not required;
10386 therefore, a shell programmer should consider a pipeline to be in a subshell environment, but
10387 not depend on it.10388 In early proposals, the description of execution environment failed to mention that each
10389 command in a multiple command pipeline could be in a subshell execution environment. For
10390 compatibility with some historical shells, the wording was phrased to allow an implementation
10391 to place any or all commands of a pipeline in the current environment. However, this means that
10392 a POSIX application must assume each command is in a subshell environment, but not depend
10393 on it.10394 The wording about shell scripts is meant to convey the fact that describing “trap actions” can
10395 only be understood in the context of the shell command language. Outside of this context, such
10396 as in a C-language program, signals are the operative condition, not traps.10397 **C.2.13 Pattern Matching Notation**10398 Pattern matching is a simpler concept and has a simpler syntax than REs, as the former is
10399 generally used for the manipulation of filenames, which are relatively simple collections of
10400 characters, while the latter is generally used to manipulate arbitrary text strings of potentially
10401 greater complexity. However, some of the basic concepts are the same, so this section points
10402 liberally to the detailed descriptions in the Base Definitions volume of IEEE Std 1003.1-2001,
10403 Chapter 9, Regular Expressions.10404 *C.2.13.1 Patterns Matching a Single Character*10405 Both quoting and escaping are described here because pattern matching must work in three
10406 separate circumstances:10407 1. Calling directly upon the shell, such as in pathname expansion or in a **case** statement. All
10408 of the following match the string or file **abc**:10409 `abc "abc" a"b" c a\bc a[b]c a["b"]c a[\b]c a["\b"]c a?c a*c`

10410 The following do not:

10411 `"a?c" a*c a\[b]c`10412 2. Calling a utility or function without going through a shell, as described for *find* and the
10413 *fnmatch()* function defined in the System Interfaces volume of IEEE Std 1003.1-2001.

10414 3. Calling utilities such as *find*, *cpio*, *tar*, or *pax* through the shell command line. In this case,
10415 shell quote removal is performed before the utility sees the argument. For example, in:

```
10416 find /bin -name "e\c[\h]o" -print
```

10417 after quote removal, the backslashes are presented to *find* and it treats them as escape
10418 characters. Both precede ordinary characters, so the *c* and *h* represent themselves and *echo*
10419 would be found on many historical systems (that have it in */bin*). To find a filename that
10420 contained shell special characters or pattern characters, both quoting and escaping are
10421 required, such as:

```
10422 pax -r ... "*a\(\?"
```

10423 to extract a filename ending with "a(?".

10424 Conforming applications are required to quote or escape the shell special characters (sometimes
10425 called metacharacters). If used without this protection, syntax errors can result or
10426 implementation extensions can be triggered. For example, the KornShell supports a series of
10427 extensions based on parentheses in patterns.

10428 The restriction on a circumflex in a bracket expression is to allow implementations that support
10429 pattern matching using the circumflex as the negation character in addition to the exclamation
10430 mark. A conforming application must use something like "[\^!]" to match either character.

10431 C.2.13.2 Patterns Matching Multiple Characters

10432 Since each asterisk matches zero or more occurrences, the patterns "a*b" and "a**b" have
10433 identical functionality.

10434 Examples

10435 a[bc] Matches the strings "ab" and "ac".

10436 a*d Matches the strings "ad", "abd", and "abcd", but not the string "abc".

10437 a*d* Matches the strings "ad", "abcd", "abcdef", "aaaad", and "adddd".

10438 *a*d Matches the strings "ad", "abcd", "efabcd", "aaaad", and "adddd".

10439 C.2.13.3 Patterns Used for Filename Expansion

10440 The caveat about a slash within a bracket expression is derived from historical practice. The
10441 pattern "a[b/c]d" does not match such pathnames as **abd** or **a/d**. On some implementations
10442 (including those conforming to the Single UNIX Specification), it matched a pathname of
10443 literally "a[b/c]d". On other systems, it produced an undefined condition (an unescaped '['
10444 used outside a bracket expression). In this version, the XSI behavior is now required.

10445 Filenames beginning with a period historically have been specially protected from view on
10446 UNIX systems. A proposal to allow an explicit period in a bracket expression to match a leading
10447 period was considered; it is allowed as an implementation extension, but a conforming
10448 application cannot make use of it. If this extension becomes popular in the future, it will be
10449 considered for a future version of the Shell and Utilities volume of IEEE Std 1003.1-2001.

10450 Historical systems have varied in their permissions requirements. To match **f*/bar** has required
10451 read permissions on the **f*** directories in the System V shell, but the Shell and Utilities volume of
10452 IEEE Std 1003.1-2001, the C shell, and KornShell require only search permissions.

10453 C.2.14 Special Built-In Utilities

10454 See the RATIONALE sections on the individual reference pages.

10455 C.3 Batch Environment Services and Utilities**10456 Scope of the Batch Environment Option**

10457 This section summarizes the deliberations of the IEEE P1003.15 (Batch Environment) working
10458 group in the development of the Batch Environment option, which covers a set of services and
10459 utilities defining a batch processing system.

10460 This informative section contains historical information concerning the contents of the
10461 amendment and describes why features were included or discarded by the working group.

10462 History of Batch Systems

10463 The supercomputing technical committee began as a “Birds Of a Feather” (BOF) at the January
10464 1987 Usenix meeting. There was enough general interest to form a supercomputing attachment
10465 to the /usr/group working groups. Several subgroups rapidly formed. Of those subgroups, the
10466 batch group was the most ambitious. The first early meetings were spent evaluating user needs
10467 and existing batch implementations.

10468 To evaluate user needs, individuals from the supercomputing community came and presented
10469 their needs. Common requests were flexibility, interoperability, control of resources, and ease-
10470 of-use. Backward-compatibility was not an issue. The working group then evaluated some
10471 existing systems. The following different systems were evaluated:

- 10472 • PROD
- 10473 • Convex Distributed Batch
- 10474 • NQS
- 10475 • CTSS
- 10476 • MDQS from Ballistics Research Laboratory (BRL)

10477 Finally, NQS was chosen as a model because it satisfied not only the most user requirements, but
10478 because it was public domain, already implemented on a variety of hardware platforms, and
10479 network-based.

10480 Historical Implementations of Batch Systems

10481 Deferred processing of work under the control of a scheduler has been a feature of most
10482 proprietary operating systems from the earliest days of multi-user systems in order to maximize
10483 utilization of the computer.

10484 The arrival of UNIX systems proved to be a dilemma to many hardware providers and users
10485 because it did not include the sophisticated batch facilities offered by the proprietary systems.
10486 This omission was rectified in 1986 by NASA Ames Research Center who developed the
10487 Network Queuing System (NQS) as a portable UNIX application that allowed the routing and
10488 processing of batch “jobs” in a network. To encourage its usage, the product was later put into
10489 the public domain. It was promptly picked up by UNIX hardware providers, and ported and
10490 developed for their respective hardware and UNIX implementations.

10491 Many major vendors, who traditionally offer a batch-dominated environment, ported the
10492 public-domain product to their systems, customized it to support the capabilities of their
10493 systems, and added many customer-requested features.

10494 Due to the strong hardware provider and customer acceptance of NQS, it was decided to use
10495 NQS as the basis for the POSIX Batch Environment amendment in 1987. Other batch systems
10496 considered at the time included CTSS, MDQS (a forerunner of NQS from the Ballistics Research
10497 Laboratory), and PROD (a Los Alamos Labs development). None were thought to have both the
10498 functionality and acceptability of NQS.

10499 **NQS Differences from the *at* utility**

10500 The base standard *at* and *batch* utilities are not sufficient to meet the batch processing needs in a
10501 supercomputing environment and additional functionality in the areas of resource management,
10502 job scheduling, system management, and control of output is required.

10503 **Batch Environment Option Definitions**

10504 The concept of a batch job is closely related to a session with a session leader. The main
10505 difference is that a batch job does not have a controlling terminal. There has been much debate
10506 over whether to use the term “request” or “job”. Job was the final choice because of the
10507 historical use of this term in the batch environment.

10508 The current definition for job identifiers is not sufficient with the model of destinations. The
10509 current definition is:

```
10510     sequence_number.originating_host
```

10511 Using the model of destination, a host may include multiple batch nodes, the location of which is
10512 identified uniquely by a name or directory service. If the current definition is used, batch nodes
10513 running on the same host would have to coordinate their use of sequence numbers, as sequence
10514 numbers are assigned by the originating host. The alternative is to use the originating batch node
10515 name instead of the originating host name.

10516 The reasons for wishing to run more than one batch system per host could be the following.

10517 A test and production batch system are maintained on a single host. This is most likely in a
10518 development facility, but could also arise when a site is moving from one version to another.
10519 The new batch system could be installed as a test version that is completely separate from the
10520 production batch system, so that problems can be isolated to the test system. Requiring the batch
10521 nodes to coordinate their use of sequence numbers creates a dependency between the two
10522 nodes, and that defeats the purpose of running two nodes.

10523 A site has multiple departments using a single host, with different management policies. An
10524 example of contention might be in job selection algorithms. One group might want a FIFO type
10525 of selection, while another group wishes to use a more complex algorithm based on resource
10526 availability. Again, requiring the batch nodes to coordinate is an unnecessary binding.

10527 The proposal eventually accepted was to replace originating host with originating batch node.
10528 This supplies sufficient granularity to ensure unique job identifiers. If more than one batch node
10529 is on a particular host, they each have their own unique name.

10530 The queue portion of a destination is not part of the job identifier as these are not required to be
10531 unique between batch nodes. For instance, two batch nodes may both have queues called small,
10532 medium, and large. It is only the batch node name that is uniquely identifiable throughout the
10533 batch system. The queue name has no additional function in this context.

10534 Assume there are three batch nodes, each of which has its own name server. On batch node one,
10535 there are no queues. On batch node two, there are fifty queues. On batch node three, there are
10536 forty queues. The system administrator for batch node one does not have to configure queues,
10537 because there are none implemented. However, if a user wishes to send a job to either batch
10538 node two or three, the system administrator for batch node one must configure a destination
10539 that maps to the appropriate batch node and queue. If every queue is to be made accessible from
10540 batch node one, the system administrator has to configure ninety destinations.

10541 To avoid requiring this, there should be a mechanism to allow a user to separate the destination
10542 into a batch node name and a queue name. Then, an implementation that is configured to get to
10543 all the batch nodes does not need any more configuration to allow a user to get to all of the
10544 queues on all of the batch nodes. The node name is used to locate the batch node, while the
10545 queue name is sent unchanged to that batch node.

10546 The following are requirements that a destination identifier must be capable of providing:

- 10547 • The ability to direct a job to a queue in a particular batch node.
- 10548 • The ability to direct a job to a particular batch node.
- 10549 • The ability to group at a higher level than just one queue. This includes grouping similar
10550 queues across multiple batch nodes (this is a pipe queue).
- 10551 • The ability to group batch nodes. This allows a user to submit a job to a group name with no
10552 knowledge of the batch node configuration. This also provides aliasing as a special case.
10553 Aliasing is a group containing only one batch node name. The group name is the alias.

10554 In addition, the administrator has the following requirements:

- 10555 • The ability to control access to the queues.
- 10556 • The ability to control access to the batch nodes.
- 10557 • The ability to control access to groups of queues (pipe queues).
- 10558 • The ability to configure retry time intervals and durations.

10559 The requirements of the user are met by destination as explained in the following.

10560 The user has the ability to specify a queue name, which is known only to the batch node
10561 specified. There is no configuration of these queues required on the submitting node.

10562 The user has the ability to specify a batch node whose name is network-unique. The
10563 configuration required is that the batch node be defined as an application, just as other
10564 applications such as FTP are configured.

10565 Once a job reaches a queue, it can again become a user of the batch system. The batch node can
10566 choose to send the job to another batch node or queue or both. In other words, the routing is at
10567 an application level, and it is up to the batch system to choose where the job will be sent.
10568 Configuration is up to the batch node where the queue resides. This provides grouping of
10569 queues across batch nodes or within a batch node. The user submits the job to a queue, which by
10570 definition routes the job to other queues or nodes or both.

10571 A node name may be given to a naming service, which returns multiple addresses as opposed to
10572 just one. This provides grouping at a batch node level. This is a local issue, meaning that the
10573 batch node must choose only one of these addresses. The list of addresses is not sent with the
10574 job, and once the job is accepted on another node, there is no connection between the list and the
10575 job. The requirements of the administrator are met by destination as explained in the following.

10576 The control of queues is a batch system issue, and will be done using the batch administrative
10577 utilities.

- 10578 The control of nodes is a network issue, and will be done through whatever network facilities
10579 are available.
- 10580 The control of access to groups of queues (pipe queues) is covered by the control of any other
10581 queue. The fact that the job may then be sent to another destination is not relevant.
- 10582 The propagation of a job across more than one point-to-point connection was dropped because
10583 of its complexity and because all of the issues arising from this capability could not be resolved.
10584 It could be provided as additional functionality at some time in the future.
- 10585 The addition of *network* as a defined term was done to clarify the difference between a network
10586 of batch nodes as opposed to a network of hosts. A network of batch nodes is referred to as a
10587 batch system. The network refers to the actual host configuration. A single host may have
10588 multiple batch nodes.
- 10589 In the absence of a standard network naming convention, this option establishes its own
10590 convention for the sake of consistency and expediency. This is subject to change, should a future
10591 working group develop a standard naming convention for network pathnames.

10592 C.3.1 Batch General Concepts

- 10593 During the development of the Batch Environment option, a number of topics were discussed at
10594 length which influenced the wording of the normative text but could not be included in the final
10595 text. The following items are some of the most significant terms and concepts of those discussed:
- 10596 • Small and Consistent Command Set

10597 Often, conventional utilities from UNIX systems have a very complicated utility syntax and
10598 usage. This can often result in confusion and errors when trying to use them. The Batch
10599 Environment option utility set, on the other hand, has been paired to a small set of robust
10600 utilities with an orthogonal calling sequence.
 - 10601 • Checkpoint/Restart

10602 This feature permits an already executing process to checkpoint or save its contents. Some
10603 implementations permit this at both the batch utility level (for example, checkpointing this
10604 job upon its abnormal termination) or from within the job itself via a system call. Support of
10605 checkpoint/restart is optional. A conscious, careful effort was made to make the *qsub* utility
10606 consistently refer to checkpoint/restart as optional functionality.
 - 10607 • Rerunability

10608 When a user submits a job for batch processing, they can designate it “rerunnable” in that it
10609 will automatically resume execution from the start of the job if the machine on which it was
10610 executing crashes for some reason. The decision on whether the job will be rerun or not is
10611 entirely up to the submitter of the job and no decisions will be made within the batch system.
10612 A job that is rerunnable and has been submitted with the proper checkpoint/restart switch
10613 will first be checkpointed and execution begun from that point. Furthermore, use of the
10614 implementation-defined checkpoint/restart feature will not be defined in this context.
 - 10615 • Error Codes

10616 All utilities exit with error status zero (0) if successful, one (1) if a user error occurred, and
10617 two (2) for an internal Batch Environment option error.
 - 10618 • Level of Portability

10619 Portability is specified at both the user, operator, and administrator levels. A conforming
10620 batch implementation prevents identical functionality and behavior at all these levels.
10621 Additionally, portable batch shell scripts with embedded Batch Environment option utilities

- 10622 add an additional level of portability.
- 10623 • Resource Specification
- 10624 A small set of globally understood resources, such as memory and CPU time, is specified. All
10625 conforming batch implementations are able to process them in a manner consistent with the
10626 yet-to-be-developed resource management model. Resources not in this amendment set are
10627 ignored and passed along as part of the argument stream of the utility.
- 10628 • Queue Position
- 10629 Queue position is the place a job occupies in a queue. It is dependent on a variety of factors
10630 such as submission time and priority. Since priority may be affected by the implementation
10631 of fair share scheduling, the definition of queue position is implementation-defined.
- 10632 • Queue ID
- 10633 A numerical queue ID is an external requirement for purposes of accounting. The
10634 identification number was chosen over queue name for processing convenience.
- 10635 • Job ID
- 10636 A common notion of “jobs” is a collection of processes whose process group cannot be
10637 altered and is used for resource management and accounting. This concept is
10638 implementation-defined and, as such, has been omitted from the batch amendment.
- 10639 • Bytes *versus* Words
- 10640 Except for one case, bytes are used as the standard unit for memory size. Furthermore, the
10641 definition of a word varies from machine to machine. Therefore, bytes will be the default unit
10642 of memory size.
- 10643 • Regular Expressions
- 10644 The standard definition of regular expressions is much too broad to be used in the batch
10645 utility syntax. All that is needed is a simple concept of “all”; for example, delete all my jobs
10646 from the named queue. For this reason, regular expressions have been eliminated from the
10647 batch amendment.
- 10648 • Display Privacy
- 10649 How much data should be displayed locally through functions? Local policy dictates the
10650 amount of privacy. Library functions must be used to create and enforce local policy.
10651 Network and local *qstats* must reflect the policy of the server machine.
- 10652 • Remote Host Naming Convention
- 10653 It was decided that host names would be a maximum of 255 characters in length, with at
10654 most 15 characters being shown in displays. The 255 character limit was chosen because it is
10655 consistent with BSD. The 15-character limit was an arbitrary decision.
- 10656 • Network Administration
- 10657 Network administration is important, but is outside the scope of the batch amendment.
10658 Network administration could be done with *rsh*. However, authentication becomes two-
10659 sided.
- 10660 • Network Administration Philosophy
- 10661 Keep it simple. Centralized management should be possible. For example, Los Alamos needs
10662 a dumb set of CPUs to be managed by a central system *versus* several independently-
10663 managed systems as is the general case for the Batch Environment option.

- 10664 • Operator Utility Defaults (that is, Default Host, User, Account, and so on)
- 10665 It was decided that usability would override orthogonality and syntactic consistency.
- 10666 • The Batch System Manager and Operator Distinction
- 10667 The distinction between manager and operator is that operators can only control the flow of
- 10668 jobs. A manager can alter the batch system configuration in addition to job flow. POSIX
- 10669 makes a distinction between user and system administrator but goes no further. The
- 10670 concepts of manager and operator privileges fall under local policy. The distinction between
- 10671 manager and operator is historical in batch environments, and the Batch Environment option
- 10672 has continued that distinction.
- 10673 • The Batch System Administrator
- 10674 An administrator is equivalent to a batch system manager.

10675 C.3.2 Batch Services

- 10676 This rationale is provided as informative rather than normative text, to avoid placing
10677 requirements on implementors regarding the use of symbolic constants, but at the same time to
10678 give implementors a preferred practice for assigning values to these constants to promote
10679 interoperability.
- 10680 The *Checkpoint* and *Minimum_Cpu_Interval* attributes induce a variety of behavior depending
10681 upon their values. Some jobs cannot or should not be checkpointed. Other users will simply
10682 need to ensure job continuation across planned downtimes; for example, scheduled preventive
10683 maintenance. For users consuming expensive resources, or for jobs that run longer than the
10684 mean time between failures, however, periodic checkpointing may be essential. However,
10685 system administrators must be able to set minimum checkpoint intervals on a queue-by-queue
10686 basis to guard against, for example, naive users specifying interval values too small on
10687 memory-intensive jobs. Otherwise, system overhead would adversely affect performance.
- 10688 The use of symbolic constants, such as `NO_CHECKPOINT`, was introduced to lend a degree of
10689 formalism and portability to this option.
- 10690 Support for checkpointing is optional for servers. However, clients must provide for the `-c`
10691 option, since in a distributed environment the job may run on a server that does provide such
10692 support, even if the host of the client does not support the checkpoint feature.
- 10693 If the user does not specify the `-c` option, the default action is left unspecified by this option.
10694 Some implementations may wish to do checkpointing by default; others may wish to checkpoint
10695 only under an explicit request from the user.
- 10696 The *Priority* attribute has been made non-optional. All clients already had been required to
10697 support the `-p` option. The concept of prioritization is common in historical implementations.
10698 The default priority is left to the server to establish.
- 10699 The *Hold_Types* attribute has been modified to allow for implementation-defined hold types to
10700 be passed to a batch server.
- 10701 It was the intent of the IEEE P1003.15 working group to mandate the support for the
10702 *Resource_List* attribute in this option by referring to another amendment, specifically the
10703 IEEE P1003.1a draft standard. However, during the development of the IEEE P1003.1a draft
10704 standard this was excluded. As such this requirement has been removed from the normative
10705 text.
- 10706 The *Shell_Path* attribute has been modified to accept a list of shell paths that are associated with
10707 a host. The name of the attribute has been changed to *Shell_Path_List*.

10708 **C.3.3 Common Behavior for Batch Environment Utilities**

10709 This section was defined to meet the goal of a “Small and Consistent Command Set” for this
10710 option.

10711 **C.4 Utilities**

10712 For the utilities included in IEEE Std 1003.1-2001, see the RATIONALE sections on the individual
10713 reference pages.

10714 **Exclusion of Utilities**

10715 The set of utilities contained in IEEE Std 1003.1-2001 is drawn from the base documents, with
10716 one addition: the *c99* utility. This section contains rationale for some of the deliberations that led
10717 to this set of utilities, and why certain utilities were excluded.

10718 Many utilities were evaluated by the standard developers; more historical utilities were
10719 excluded from the base documents than included. The following list contains many common
10720 UNIX system utilities that were not included as mandatory utilities, in the User Portability
10721 Utilities option, in the XSI extension, or in one of the software development groups. It is
10722 logistically difficult for this rationale to distribute correctly the reasons for not including a utility
10723 among the various utility options. Therefore, this section covers the reasons for all utilities not
10724 included in IEEE Std 1003.1-2001.

10725 This rationale is limited to a discussion of only those utilities actively or indirectly evaluated by
10726 the standard developers of the base documents, rather than the list of all known UNIX utilities
10727 from all its variants.

10728 *adb* The intent of the various software development utilities was to assist in the
10729 installation (rather than the actual development and debugging) of applications.
10730 This utility is primarily a debugging tool. Furthermore, many useful aspects of *adb*
10731 are very hardware-specific.

10732 *as* Assemblers are hardware-specific and are included implicitly as part of the
10733 compilers in IEEE Std 1003.1-2001.

10734 *banner* The only known use of this command is as part of the *lp* printer header pages. It
10735 was decided that the format of the header is implementation-defined, so this utility
10736 is superfluous to application portability.

10737 *calendar* This reminder service program is not useful to conforming applications.

10738 *cancel* The *lp* (line printer spooling) system specified is the most basic possible and did
10739 not need this level of application control.

10740 *chroot* This is primarily of administrative use, requiring superuser privileges.

10741 *col* No utilities defined in IEEE Std 1003.1-2001 produce output requiring such a filter.
10742 The *nroff* text formatter is present on many historical systems and will continue to
10743 remain as an extension; *col* is expected to be shipped by all the systems that ship
10744 *nroff*.

10745 *cpio* This has been replaced by *pax*, for reasons explained in the rationale for that utility.

10746 *cpp* This is subsumed by *c99*.

10747 *cu* This utility is terminal-oriented and is not useful from shell scripts or typical
10748 application programs.

10749	<i>dc</i>	The functionality of this utility can be provided by the <i>bc</i> utility; <i>bc</i> was selected because it was easier to use and had superior functionality. Although the historical versions of <i>bc</i> are implemented using <i>dc</i> as a base, IEEE Std 1003.1-2001 prescribes the interface and not the underlying mechanism used to implement it.
10750		
10751		
10752		
10753	<i>dircmp</i>	Although a useful concept, the historical output of this directory comparison program is not suitable for processing in application programs. Also, the <i>diff -r</i> command gives equivalent functionality.
10754		
10755		
10756	<i>dis</i>	Disassemblers are hardware-specific.
10757	<i>emacs</i>	The community of <i>emacs</i> editing enthusiasts was adamant that the full <i>emacs</i> editor not be included in the base documents because they were concerned that an attempt to standardize this very powerful environment would encourage vendors to ship versions conforming strictly to the standard, but lacking the extensibility required by the community. The author of the original <i>emacs</i> program also expressed his desire to omit the program. Furthermore, there were a number of historical UNIX systems that did not include <i>emacs</i> , or included it without supporting it, but there were very few that did not include and support <i>vi</i> .
10758		
10759		
10760		
10761		
10762		
10763		
10764		
10765	<i>ld</i>	This is subsumed by <i>c99</i> .
10766	<i>line</i>	The functionality of <i>line</i> can be provided with <i>read</i> .
10767	<i>lint</i>	This technology is partially subsumed by <i>c99</i> . It is also hard to specify the degree of checking for possible error conditions in programs in any compiler, and specifying what <i>lint</i> would do in these cases is equally difficult.
10768		
10769		It is fairly easy to specify what a compiler does. It requires specifying the language, what it does with that language, and stating that the interpretation of any incorrect program is unspecified. Unfortunately, any description of <i>lint</i> is required to specify what to do with erroneous programs. Since the number of possible errors and questionable programming practices is infinite, one cannot require <i>lint</i> to detect all errors of any given class.
10770		
10771		
10772		
10773		
10774		Additionally, some vendors complained that since many compilers are distributed in a binary form without a <i>lint</i> facility (because the ISO C standard does not require one), implementing the standard as a stand-alone product will be much harder. Rather than being able to build upon a standard compiler component (simply by providing <i>c99</i> as an interface), source to that compiler would most likely need to be modified to provide the <i>lint</i> functionality. This was considered a major burden on system providers for a very small gain to developers (users).
10776		
10777		
10778		
10779		
10780		This utility is terminal-oriented and is not useful from shell scripts or typical application programs.
10781	<i>login</i>	
10782		This utility is an aid in creating an implementation-defined detail of object libraries that the standard developers did not feel required standardization.
10783	<i>lorder</i>	
10784		The <i>lp</i> system specified is the most basic possible and did not need this level of application control.
10785	<i>lpstat</i>	
10786		This utility was omitted in favor of <i>mailx</i> because there was a considerable functionality overlap between the two.
10787	<i>mail</i>	
10788		This was omitted in favor of <i>mkfifo</i> , as <i>mknod</i> has too many implementation-defined functions.
10789	<i>mknod</i>	
10790		
10791		
10792		

10793	<i>news</i>	This utility is terminal-oriented and is not useful from shell scripts or typical application programs.
10794		
10795	<i>pack</i>	This compression program was considered inferior to <i>compress</i> .
10796	<i>passwd</i>	This utility was proposed in a historical draft of the base documents but met with too many objections to be included. There were various reasons:
10797		
10798		<ul style="list-style-type: none"> • Changing a password should not be viewed as a command, but as part of the login sequence. Changing a password should only be done while a trusted path is in effect.
10799		
10800		
10801		<ul style="list-style-type: none"> • Even though the text in early drafts was intended to allow a variety of implementations to conform, the security policy for one site may differ from another site running with identical hardware and software. One site might use password authentication while the other did not. Vendors could not supply a <i>passwd</i> utility that would conform to IEEE Std 1003.1-2001 for all sites using their system.
10802		
10803		
10804		
10805		
10806		
10807		<ul style="list-style-type: none"> • This is really a subject for a system administration working group or a security working group.
10808		
10809	<i>pcat</i>	This compression program was considered inferior to <i>zcat</i> .
10810	<i>pg</i>	This duplicated many of the features of the <i>more</i> pager, which was preferred by the standard developers.
10811		
10812	<i>prof</i>	The intent of the various software development utilities was to assist in the installation (rather than the actual development and debugging) of applications. This utility is primarily a debugging tool.
10813		
10814		
10815	RCS	RCS was originally considered as part of a version control utilities portion of the scope. However, this aspect was abandoned by the standard developers. SCCS is now included as an optional part of the XSI extension.
10816		
10817		
10818	<i>red</i>	Restricted editor. This was not considered by the standard developers because it never provided the level of security restriction required.
10819		
10820	<i>rsh</i>	Restricted shell. This was not considered by the standard developers because it does not provide the level of security restriction that is implied by historical documentation.
10821		
10822		
10823	<i>sdb</i>	The intent of the various software development utilities was to assist in the installation (rather than the actual development and debugging) of applications. This utility is primarily a debugging tool. Furthermore, some useful aspects of <i>sdb</i> are very hardware-specific.
10824		
10825		
10826		
10827	<i>sdiff</i>	The “side-by-side <i>diff</i> ” utility from System V was omitted because it is used infrequently, and even less so by conforming applications. Despite being in System V, it is not in the SVID or XPG.
10828		
10829		
10830	<i>shar</i>	Any of the numerous “shell archivers” were excluded because they did not meet the requirement of existing practice.
10831		
10832	<i>shl</i>	This utility is terminal-oriented and is not useful from shell scripts or typical application programs. The job control aspects of the shell command language are generally more useful.
10833		
10834		
10835	<i>size</i>	The intent of the various software development utilities was to assist in the installation (rather than the actual development and debugging) of applications.
10836		

10837		This utility is primarily a debugging tool.
10838	<i>spell</i>	This utility is not useful from shell scripts or typical application programs. The <i>spell</i> utility was considered, but was omitted because there is no known technology that can be used to make it recognize general language for user-specified input without providing a complete dictionary along with the input file.
10839		
10840		
10841		
10842	<i>su</i>	This utility is not useful from shell scripts or typical application programs. (There was also sentiment to avoid security-related utilities.)
10843		
10844	<i>sum</i>	This utility was renamed <i>cksum</i> .
10845	<i>tar</i>	This has been replaced by <i>pax</i> , for reasons explained in the rationale for that utility.
10846	<i>tsort</i>	This utility is an aid in creating an implementation-defined detail of object libraries that the standard developers did not feel required standardization.
10847		
10848	<i>unpack</i>	This compression program was considered inferior to <i>uncompress</i> .
10849	<i>wall</i>	This utility is terminal-oriented and is not useful in shell scripts or typical applications. It is generally used only by system administrators.
10850		

10851 / *Rationale (Informative)*

10852 **Part D:**

10853 **Portability Considerations**

10854 *The Open Group*

10855 *The Institute of Electrical and Electronics Engineers, Inc.*

Portability Considerations (Informative)

10856

10857 This section contains information to satisfy various international requirements:

- 10858 • Section D.1 describes perceived user requirements.
- 10859 • Section D.2 (on page 270) indicates how the facilities of IEEE Std 1003.1-2001 satisfy those
10860 requirements.
- 10861 • Section D.3 (on page 277) offers guidance to writers of profiles on how the configurable
10862 options, limits, and optional behavior of IEEE Std 1003.1-2001 should be cited in profiles.

10863 D.1 User Requirements

10864 This section describes the user requirements that were perceived by the developers of
10865 IEEE Std 1003.1-2001. The primary source for these requirements was an analysis of historical
10866 practice in widespread use, as typified by the base documents listed in Section A.1.1 (on page 3).

10867 IEEE Std 1003.1-2001 addresses the needs of users requiring open systems solutions for source
10868 code portability of applications. It currently addresses users requiring open systems solutions
10869 for source-code portability of applications involving multi-programming and process
10870 management (creating processes, signaling, and so on); access to files and directories in a
10871 hierarchy of file systems (opening, reading, writing, deleting files, and so on); access to
10872 asynchronous communications ports and other special devices; access to information about
10873 other users of the system; facilities supporting applications requiring bounded (realtime)
10874 response.

10875 The following users are identified for IEEE Std 1003.1-2001:

- 10876 • Those employing applications written in high-level languages, such as C, Ada, or FORTRAN.
- 10877 • Users who desire conforming applications that do not necessarily require the characteristics
10878 of high-level languages (for example, the speed of execution of compiled languages or the
10879 relative security of source code intellectual property inherent in the compilation process).
- 10880 • Users who desire conforming applications that can be developed quickly and can be
10881 modified readily without the use of compilers and other system components that may be
10882 unavailable on small systems or those without special application development capabilities.
- 10883 • Users who interact with a system to achieve general-purpose time-sharing capabilities
10884 common to most business or government offices or academic environments: editing, filing,
10885 inter-user communications, printing, and so on.
- 10886 • Users who develop applications for POSIX-conformant systems.
- 10887 • Users who develop applications for UNIX systems.

10888 An acknowledged restriction on applicable users is that they are limited to the group of
10889 individuals who are familiar with the style of interaction characteristic of historically-derived
10890 systems based on one of the UNIX operating systems (as opposed to other historical systems
10891 with different models, such as MS/DOS, Macintosh, VMS, MVS, and so on). Typical users
10892 would include program developers, engineers, or general-purpose time-sharing users.

10893 The requirements of users of IEEE Std 1003.1-2001 can be summarized as a single goal:
10894 *application source portability*. The requirements of the user are stated in terms of the requirements

10895 of portability of applications. This in turn becomes a requirement for a standardized set of
10896 syntax and semantics for operations commonly found on many operating systems.

10897 The following sections list the perceived requirements for application portability.

10898 **D.1.1 Configuration Interrogation**

10899 An application must be able to determine whether and how certain optional features are
10900 provided and to identify the system upon which it is running, so that it may appropriately adapt
10901 to its environment.

10902 Applications must have sufficient information to adapt to varying behaviors of the system.

10903 **D.1.2 Process Management**

10904 An application must be able to manage itself, either as a single process or as multiple processes.
10905 Applications must be able to manage other processes when appropriate.

10906 Applications must be able to identify, control, create, and delete processes, and there must be
10907 communication of information between processes and to and from the system.

10908 Applications must be able to use multiple flows of control with a process (threads) and
10909 synchronize operations between these flows of control.

10910 **D.1.3 Access to Data**

10911 Applications must be able to operate on the data stored on the system, access it, and transmit it
10912 to other applications. Information must have protection from unauthorized or accidental access
10913 or modification.

10914 **D.1.4 Access to the Environment**

10915 Applications must be able to access the external environment to communicate their input and
10916 results.

10917 **D.1.5 Access to Determinism and Performance Enhancements**

10918 Applications must have sufficient control of resource allocation to ensure the timeliness of
10919 interactions with external objects.

10920 **D.1.6 Operating System-Dependent Profile**

10921 The capabilities of the operating system may make certain optional characteristics of the base
10922 language in effect no longer optional, and this should be specified.

10923 **D.1.7 I/O Interaction**

10924 The interaction between the C language I/O subsystem (*stdio*) and the I/O subsystem of
10925 IEEE Std 1003.1-2001 must be specified.

10926 D.1.8 Internationalization Interaction

10927 The effects of the environment of IEEE Std 1003.1-2001 on the internationalization facilities of the
10928 C language must be specified.

10929 D.1.9 C-Language Extensions

10930 Certain functions in the C language must be extended to support the additional capabilities
10931 provided by IEEE Std 1003.1-2001.

10932 D.1.10 Command Language

10933 Users should be able to define procedures that combine simple tools and/or applications into
10934 higher-level components that perform to the specific needs of the user. The user should be able
10935 to store, recall, use, and modify these procedures. These procedures should employ a powerful
10936 command language that is used for recurring tasks in conforming applications (scripts) in the
10937 same way that it is used interactively to accomplish one-time tasks. The language and the
10938 utilities that it uses must be consistent between systems to reduce errors and retraining.

10939 D.1.11 Interactive Facilities

10940 Use the system to accomplish individual tasks at an interactive terminal. The interface should be
10941 consistent, intuitive, and offer usability enhancements to increase the productivity of terminal
10942 users, reduce errors, and minimize retraining costs. Online documentation or usage assistance
10943 should be available.

10944 D.1.12 Accomplish Multiple Tasks Simultaneously

10945 Access applications and interactive facilities from a single terminal without requiring serial
10946 execution: switch between multiple interactive tasks; schedule one-time or periodic background
10947 work; display the status of all work in progress or scheduled; influence the priority scheduling of
10948 work, when authorized.

10949 D.1.13 Complex Data Manipulation

10950 Manipulate data in files in complex ways: sort, merge, compare, translate, edit, format, pattern
10951 match, select subsets (strings, columns, fields, rows, and so on). These facilities should be
10952 available to both conforming applications and interactive users.

10953 D.1.14 File Hierarchy Manipulation

10954 Create, delete, move/rename, copy, backup/archive, and display files and directories. These
10955 facilities should be available to both conforming applications and interactive users.

10956 D.1.15 Locale Configuration

10957 Customize applications and interactive sessions for the cultural and language conventions of the
10958 user. Employ a wide variety of standard character encodings. These facilities should be available
10959 to both conforming applications and interactive users.

10960 D.1.16 Inter-User Communication

10961 Send messages or transfer files to other users on the same system or other systems on a network.
10962 These facilities should be available to both conforming applications and interactive users.

10963 D.1.17 System Environment

10964 Display information about the status of the system (activities of users and their interactive and
10965 background work, file system utilization, system time, configuration, and presence of optional
10966 facilities) and the environment of the user (terminal characteristics, and so on). Inform the
10967 system operator/administrator of problems. Control access to user files and other resources.

10968 D.1.18 Printing

10969 Output files on a variety of output device classes, accessing devices on local or network-
10970 connected systems. Control (or influence) the formatting, priority scheduling, and output
10971 distribution of work. These facilities should be available to both conforming applications and
10972 interactive users.

10973 D.1.19 Software Development

10974 Develop (create and manage source files, compile/interpret, debug) portable open systems
10975 applications and package them for distribution to, and updating of, other systems.

10976 D.2 Portability Capabilities

10977 This section describes the significant portability capabilities of IEEE Std 1003.1-2001 and
10978 indicates how the user requirements listed in Section D.1 (on page 267) are addressed. The
10979 capabilities are listed in the same format as the preceding user requirements; they are
10980 summarized below:

- 10981 • Configuration Interrogation
- 10982 • Process Management
- 10983 • Access to Data
- 10984 • Access to the Environment
- 10985 • Access to Determinism and Performance Enhancements
- 10986 • Operating System-Dependent Profile
- 10987 • I/O Interaction
- 10988 • Internationalization Interaction
- 10989 • C-Language Extensions
- 10990 • Command Language
- 10991 • Interactive Facilities
- 10992 • Accomplish Multiple Tasks Simultaneously
- 10993 • Complex Data Manipulation
- 10994 • File Hierarchy Manipulation

- 10995 • Locale Configuration
- 10996 • Inter-User Communication
- 10997 • System Environment
- 10998 • Printing
- 10999 • Software Development

11000 D.2.1 Configuration Interrogation

11001 The *uname()* operation provides basic identification of the system. The *sysconf()*, *pathconf()*, and
 11002 *fpathconf()* functions and the *getconf* utility provide means to interrogate the implementation to
 11003 determine how to adapt to the environment in which it is running. These values can be either
 11004 static (indicating that all instances of the implementation have the same value) or dynamic
 11005 (indicating that different instances of the implementation have the different values, or that the
 11006 value may vary for other reasons, such as reconfiguration).

11007 Unsatisfied Requirements

11008 None directly. However, as new areas are added, there will be a need for additional capability in
 11009 this area.

11010 D.2.2 Process Management

11011 The *fork()*, *exec* family, *posix_spawn()*, and *posix_spawnp()* functions provide for the creation of
 11012 new processes or the insertion of new applications into existing processes. The *_Exit()*, *_exit()*,
 11013 *exit()*, and *abort()* functions allow for the termination of a process by itself. The *wait()* and
 11014 *waitpid()* functions allow one process to deal with the termination of another.

11015 The *times()* function allows for basic measurement of times used by a process. Various
 11016 functions, including *fstat()*, *getegid()*, *geteuid()*, *getgid()*, *getgrgid()*, *getgrnam()*, *getlogin()*,
 11017 *getpid()*, *getppid()*, *getpwnam()*, *getpwuid()*, *getuid()*, *lstat()*, and *stat()*, provide for access to the
 11018 identifiers of processes and the identifiers and names of owners of processes (and files).

11019 The various functions operating on environment variables provide for communication of
 11020 information (primarily user-configurable defaults) from a parent to child processes.

11021 The operations on the current working directory control and interrogate the directory from
 11022 which relative filename searches start. The *umask()* function controls the default protections
 11023 applied to files created by the process.

11024 The *alarm()*, *pause()*, *sleep()*, *ualarm()*, and *usleep()* operations allow the process to suspend until
 11025 a timer has expired or to be notified when a period of time has elapsed. The *time()* operation
 11026 interrogates the current time and date.

11027 The signal mechanism provides for communication of events either from other processes or
 11028 from the environment to the application, and the means for the application to control the effect
 11029 of these events. The mechanism provides for external termination of a process and for a process
 11030 to suspend until an event occurs. The mechanism also provides for a value to be associated with
 11031 an event.

11032 Job control provides a means to group processes and control them as groups, and to control their
 11033 access to the function between the user and the system (the “controlling terminal”). It also
 11034 provides the means to suspend and resume processes.

11035 The Process Scheduling option provides control of the scheduling and priority of a process.

11036 The Message Passing option provides a means for interprocess communication involving small
11037 amounts of data.

11038 The Memory Management facilities provide control of memory resources and for the sharing of
11039 memory. This functionality is mandatory on XSI-conformant systems.

11040 The Threads facilities provide multiple flows of control with a process (threads),
11041 synchronization between threads, association of data with threads, and controlled cancelation of
11042 threads.

11043 The XSI interprocess communications functionality provide an alternate set of facilities to
11044 manipulate semaphores, message queues, and shared memory. These are provided on XSI-
11045 conformant systems to support conforming applications developed to run on UNIX systems.

11046 **D.2.3 Access to Data**

11047 The *open()*, *close()*, *fclose()*, *fopen()*, and *pipe()* functions provide for access to files and data.
11048 Such files may be regular files, interprocess data channels (pipes), or devices. Additional types
11049 of objects in the file system are permitted and are being contemplated for standardization.

11050 The *access()*, *chmod()*, *chown()*, *dup()*, *dup2()*, *fchmod()*, *fcntl()*, *fstat()*, *ftruncate()*, *lstat()*,
11051 *readlink()*, *realpath()*, *stat()*, and *utime()* functions allow for control and interrogation of file and
11052 file-related objects (including symbolic links), and their ownership, protections, and timestamps.

11053 The *fgetc()*, *fputc()*, *fread()*, *fseek()*, *fsetpos()*, *fwrite()*, *getc()*, *getchar()*, *lseek()*, *putchar()*, *putc()*,
11054 *read()*, and *write()* functions provide for data transfer from the application to files (in all their
11055 forms).

11056 The *closedir()*, *link()*, *mkdir()*, *opendir()*, *readdir()*, *rename()*, *rmdir()*, *rewinddir()*, and *unlink()*
11057 functions provide for a complete set of operations on directories. Directories can arbitrarily
11058 contain other directories, and a single file can be mentioned in more than one directory.

11059 The file-locking mechanism provides for advisory locking (protection during transactions) of
11060 ranges of bytes (in effect, records) in a file.

11061 The *confstr()*, *fpathconf()*, *pathconf()*, and *sysconf()* functions provide for enquiry as to the
11062 behavior of the system where variability is permitted.

11063 The Synchronized Input and Output option provides for assured commitment of data to media.

11064 The Asynchronous Input and Output option provides for initiation and control of asynchronous
11065 data transfers.

11066 **D.2.4 Access to the Environment**

11067 The operations and types in the Base Definitions volume of IEEE Std 1003.1-2001, Chapter 11,
11068 General Terminal Interface are provided for access to asynchronous serial devices. The primary
11069 intended use for these is the controlling terminal for the application (the interaction point
11070 between the user and the system). They are general enough to be used to control any
11071 asynchronous serial device. The functions are also general enough to be used with many other
11072 device types as a user interface when some emulation is provided.

11073 Less detailed access is provided for other device types, but in many instances an application
11074 need not know whether an object in the file system is a device or a regular file to operate
11075 correctly.

11076 **Unsatisfied Requirements**

11077 Detailed control of common device classes, specifically magnetic tape, is not provided.

11078 **D.2.5 Bounded (Realtime) Response**

11079 The Realtime Signals Extension provides queued signals and the prioritization of the handling of
 11080 signals. The SCHED_FIFO, SCHED_SPORADIC, and SCHED_RR scheduling policies provide
 11081 control over processor allocation. The Semaphores option provides high-performance
 11082 synchronization. The Memory Management functions provide memory locking for control of
 11083 memory allocation, file mapping for high-performance, and shared memory for high-
 11084 performance interprocess communication. The Message Passing option provides for interprocess
 11085 communication without being dependent on shared memory.

11086 The Timers option provides a high resolution function called *nanosleep()* with a finer resolution
 11087 than the *sleep()* function.

11088 The Typed Memory Objects option, the Monotonic Clock option, and the Timeouts option
 11089 provide further facilities for applications to use to obtain predictable bounded response.

11090 **D.2.6 Operating System-Dependent Profile**

11091 IEEE Std 1003.1-2001 makes no distinction between text and binary files. The values of
 11092 EXIT_SUCCESS and EXIT_FAILURE are further defined.

11093 **Unsatisfied Requirements**

11094 None known, but the ISO C standard may contain some additional options that could be
 11095 specified.

11096 **D.2.7 I/O Interaction**

11097 IEEE Std 1003.1-2001 defines how each of the ISO C standard *stdio* functions interact with the
 11098 POSIX.1 operations, typically specifying the behavior in terms of POSIX.1 operations.

11099 **Unsatisfied Requirements**

11100 None.

11101 **D.2.8 Internationalization Interaction**

11102 The IEEE Std 1003.1-2001 environment operations provide a means to define the environment
 11103 for *setlocale()* and time functions such as *ctime()*. The *tzset()* function is provided to set time
 11104 conversion information.

11105 The *nl_langinfo()* function is provided as an XSI extension to query locale-specific cultural
 11106 settings.

11107 **Unsatisfied Requirements**

11108 None.

11109 D.2.9 C-Language Extensions

11110 The *setjmp()* and *longjmp()* functions are not defined to be cognizant of the signal masks defined
11111 for POSIX.1. The *sigsetjmp()* and *siglongjmp()* functions are provided to fill this gap.

11112 The *_setjmp()* and *_longjmp()* functions are provided as XSI extensions to support historic
11113 practice.

11114 Unsatisfied Requirements

11115 None.

11116 D.2.10 Command Language

11117 The shell command language, as described in the Shell and Utilities volume of
11118 IEEE Std 1003.1-2001, Chapter 2, Shell Command Language, is a common language useful in
11119 batch scripts, through an API to high-level languages (for the C-Language Binding option,
11120 *system()* and *popen()*) and through an interactive terminal (see the *sh* utility). The shell language
11121 has many of the characteristics of a high-level language, but it has been designed to be more
11122 suitable for user terminal entry and includes interactive debugging facilities. Through the use of
11123 pipelining, many complex commands can be constructed from combinations of data filters and
11124 other common components. Shell scripts can be created, stored, recalled, and modified by the
11125 user with simple editors.

11126 In addition to the basic shell language, the following utilities offer features that simplify and
11127 enhance programmatic access to the utilities and provide features normally found only in high-
11128 level languages: *basename*, *bc*, *command*, *dirname*, *echo*, *env*, *expr*, *false*, *printf*, *read*, *sleep*, *tee*, *test*,
11129 *time**,² *true*, *wait*, *xargs*, and all of the special built-in utilities in the Shell and Utilities volume of
11130 IEEE Std 1003.1-2001, Section 2.14, Special Built-In Utilities.

11131 Unsatisfied Requirements

11132 None.

11133 D.2.11 Interactive Facilities

11134 The utilities offer a common style of command-line interface through conformance to the Utility
11135 Syntax Guidelines (see the Base Definitions volume of IEEE Std 1003.1-2001, Section 12.2, Utility
11136 Syntax Guidelines) and the common utility defaults (see the Shell and Utilities volume of
11137 IEEE Std 1003.1-2001, Section 1.11, Utility Description Defaults). The *sh* utility offers an
11138 interactive command-line history and editing facility. The following utilities in the User
11139 Portability Utilities option have been customized for interactive use: *alias*, *ex*, *fc*, *mailx*, *more*, *talk*,
11140 *vi*, *unalias*, and *write*; the *man* utility offers online access to system documentation.

11141 _____

11142 2. The utilities listed with an asterisk here and later in this section are present only on systems which support the User Portability
11143 Utilities option. There may be further restrictions on the utilities offered with various configuration option combinations; see the
11144 individual utility descriptions.

11145 **Unsatisfied Requirements**

11146 The command line interface to individual utilities is as intuitive and consistent as historical
 11147 practice allows. Work underway based on graphical user interfaces may be more suitable for
 11148 novice or occasional users of the system.

11149 **D.2.12 Accomplish Multiple Tasks Simultaneously**

11150 The shell command language offers background processing through the asynchronous list
 11151 command form; see the Shell and Utilities volume of IEEE Std 1003.1-2001, Section 2.9, Shell
 11152 Commands. The *nohup* utility makes background processing more robust and usable. The *kill*
 11153 utility can terminate background jobs. When the User Portability Utilities option is supported,
 11154 the following utilities allow manipulation of jobs: *bg*, *fg*, and *jobs*. Also, if the User Portability
 11155 Utilities option is supported, the following can support periodic job scheduling, control, and
 11156 display: *at*, *batch*, *crontab*, *nice*, *ps*, and *renice*.

11157 **Unsatisfied Requirements**

11158 Terminals with multiple windows may be more suitable for some multi-tasking interactive uses
 11159 than the job control approach in IEEE Std 1003.1-2001. See the comments on graphical user
 11160 interfaces in Section D.2.11 (on page 274). The *nice* and *renice* utilities do not necessarily take
 11161 advantage of complex system scheduling algorithms that are supported by the realtime options
 11162 within IEEE Std 1003.1-2001.

11163 **D.2.13 Complex Data Manipulation**

11164 The following utilities address user requirements in this area: *asa*, *awk*, *bc*, *cmp*, *comm*, *csplit**, *cut*,
 11165 *dd*, *diff*, *ed*, *ex**, *expand**, *expr*, *find*, *fold*, *grep*, *head*, *join*, *od*, *paste*, *pr*, *printf*, *sed*, *sort*, *split**, *tabs**, *tail*,
 11166 *tr*, *unexpand**, *uniq*, *uudecode**, *uuencode**, and *wc*.

11167 **Unsatisfied Requirements**

11168 Sophisticated text formatting utilities, such as *troff* or *TeX*, are not included. Standards work in
 11169 the area of SGML may satisfy this.

11170 **D.2.14 File Hierarchy Manipulation**

11171 The following utilities address user requirements in this area: *basename*, *cd*, *chgrp*, *chmod*, *chown*,
 11172 *cksum*, *cp*, *dd*, *df**, *diff*, *dirname*, *du**, *find*, *ls*, *ln*, *mkdir*, *mkfifo*, *mv*, *patch**, *pathchk*, *pax*, *pwd*, *rm*, *rmdir*,
 11173 *test*, and *touch*.

11174 **Unsatisfied Requirements**

11175 Some graphical user interfaces offer more intuitive file manager components that allow file
 11176 manipulation through the use of icons for novice users.

11177 D.2.15 Locale Configuration

11178 The standard utilities are affected by the various *LC_* variables to achieve locale-dependent
 11179 operation: character classification, collation sequences, regular expressions and shell pattern
 11180 matching, date and time formats, numeric formatting, and monetary formatting. When the
 11181 POSIX2_LOCALEDEF option is supported, applications can provide their own locale definition
 11182 files. The following utilities address user requirements in this area: *date*, *ed*, *ex**, *find*, *grep*, *locale*,
 11183 *localedef*, *more**, *sed*, *sh*, *sort*, *tr*, *uniq*, and *vi**.

11184 The *iconv()*, *iconv_close()*, and *iconv_open()* functions are available to allow an application to
 11185 convert character data between supported character sets.

11186 The *genccat* utility and the *catopen()*, *catclose()*, and *catgets()* functions for message catalog
 11187 manipulation are available on XSI-conformant systems.

11188 Unsatisfied Requirements

11189 Some aspects of multi-byte character and state-encoded character encodings have not yet been
 11190 addressed. The C-language functions, such as *getopt()*, are generally limited to single-byte
 11191 characters. The effect of the *LC_MESSAGES* variable on message formats is only suggested at
 11192 this time.

11193 D.2.16 Inter-User Communication

11194 The following utilities address user requirements in this area: *cksum*, *mailx**, *mesg**, *patch**, *pax*,
 11195 *talk**, *uudecode**, *uuencode**, *who**, and *write**.

11196 The historical UUCP utilities are included on XSI-conformant systems.

11197 Unsatisfied Requirements

11198 None.

11199 D.2.17 System Environment

11200 The following utilities address user requirements in this area: *chgrp*, *chmod*, *chown*, *df**, *du**, *env*,
 11201 *getconf*, *id*, *logger*, *logname*, *mesg**, *newgrp**, *ps**, *stty*, *tput**, *tty*, *umask*, *uname*, and *who**.

11202 The *closelog()*, *openlog()*, *setlogmask()*, and *syslog()* functions provide System Logging facilities
 11203 on XSI-conformant systems; these are analogous to the *logger* utility.

11204 Unsatisfied Requirements

11205 None.

11206 D.2.18 Printing

11207 The following utilities address user requirements in this area: *pr* and *lp*.

11208 Unsatisfied Requirements

11209 There are no features to control the formatting or scheduling of the print jobs.

11210 D.2.19 Software Development

11211 The following utilities address user requirements in this area: *ar*, *asa*, *awk*, *c99*, *ctags**, *fort77*,
11212 *getconf*, *getopts*, *lex*, *localedef*, *make*, *nm**, *od*, *patch**, *pax*, *strings**, *strip*, *time**, and *yacc*.

11213 The *system()*, *popen()*, *pclose()*, *regcomp()*, *regex()*, *regerror()*, *regfree()*, *fnmatch()*, *getopt()*,
11214 *glob()*, *globfree()*, *wordexp()*, and *wordfree()* functions allow C-language programmers to access
11215 some of the interfaces used by the utilities, such as argument processing, regular expressions,
11216 and pattern matching.

11217 The SCCS source-code control system utilities are available on systems supporting the XSI
11218 Development option.

11219 Unsatisfied Requirements

11220 There are no language-specific development tools related to languages other than C and
11221 FORTRAN. The C tools are more complete and varied than the FORTRAN tools. There is no
11222 data dictionary or other CASE-like development tools.

11223 D.2.20 Future Growth

11224 It is arguable whether or not all functionality to support applications is potentially within the
11225 scope of IEEE Std 1003.1-2001. As a simple matter of practicality, it cannot be. Areas such as
11226 graphics, application domain-specific functionality, windowing, and so on, should be in unique
11227 standards. As such, they are properly “Unsatisfied Requirements” in terms of providing fully
11228 conforming applications, but ones which are outside the scope of IEEE Std 1003.1-2001.

11229 However, as the standards evolve, certain functionality once considered “exotic” enough to be
11230 part of a separate standard become common enough to be included in a core standard such as
11231 this. Realtime and networking, for example, have both moved from separate standards (with
11232 much difficult cross-referencing) into IEEE Std 1003.1 over time, and although no specific areas
11233 have been identified for inclusion in future revisions, such inclusions seem likely.

11234 D.3 Profiling Considerations

11235 This section offers guidance to writers of profiles on how the configurable options, limits, and
11236 optional behavior of IEEE Std 1003.1-2001 should be cited in profiles. Profile writers should
11237 consult the general guidance in POSIX.0 when writing POSIX Standardized Profiles.

11238 The information in this section is an inclusive list of features that should be considered by profile
11239 writers. Subsetting of IEEE Std 1003.1-2001 should follow the Base Definitions volume of
11240 IEEE Std 1003.1-2001, Section 2.1.5.1, Subprofiling Considerations. A set of profiling options is
11241 described in Appendix E (on page 291).

11242 D.3.1 Configuration Options

11243 There are two set of options suggested by IEEE Std 1003.1-2001: those for POSIX-conforming
11244 systems and those for X/Open System Interface (XSI) conformance. The requirements for XSI
11245 conformance are documented in the Base Definitions volume of IEEE Std 1003.1-2001 and not
11246 discussed further here, as they superset the POSIX conformance requirements.

11247 D.3.2 Configuration Options (Shell and Utilities)

11248 There are three broad optional configurations for the Shell and Utilities volume of
11249 IEEE Std 1003.1-2001: basic execution system, development system, and user portability
11250 interactive system. The options to support these, and other minor configuration options, are
11251 listed in the Base Definitions volume of IEEE Std 1003.1-2001, Chapter 2, Conformance. Profile
11252 writers should consult the following list and the comments concerning user requirements
11253 addressed by various components in Section D.2 (on page 270).

11254 POSIX2_UPE

11255 The system supports the User Portability Utilities option.

11256 This option is a requirement for a user portability interactive system. It is required
11257 frequently except for those systems, such as embedded realtime or dedicated application
11258 systems, that support little or no interactive time-sharing work by users or operators. XSI-
11259 conformant systems support this option.

11260 POSIX2_SW_DEV

11261 The system supports the Software Development Utilities option.

11262 This option is required by many systems, even those in which actual software development
11263 does not occur. The *make* utility, in particular, is required by many application software
11264 packages as they are installed onto the system. If POSIX2_C_DEV is supported,
11265 POSIX2_SW_DEV is almost a mandatory requirement because of *ar* and *make*.

11266 POSIX2_C_BIND

11267 The system supports the C-Language Bindings option.

11268 This option is required on some implementations developing complex C applications or on
11269 any system installing C applications in source form that require the functions in this option.
11270 The *system()* and *popen()* functions, in particular, are widely used by applications; the
11271 others are rather more specialized.

11272 POSIX2_C_DEV

11273 The system supports the C-Language Development Utilities option.

11274 This option is required by many systems, even those in which actual C-language software
11275 development does not occur. The *c99* utility, in particular, is required by many application
11276 software packages as they are installed onto the system. The *lex* and *yacc* utilities are used
11277 less frequently.

11278 POSIX2_FORT_DEV

11279 The system supports the FORTRAN Development Utilities option

11280 As with C, this option is needed on any system developing or installing FORTRAN
11281 applications in source form.

11282 POSIX2_FORT_RUN

11283 The system supports the FORTRAN Runtime Utilities option.

11284 This option is required for some FORTRAN applications that need the *asa* utility to convert
11285 Hollerith printing statement output. It is unknown how frequently this occurs.

11286 POSIX2_LOCALEDEF

11287 The system supports the creation of locales.

11288 This option is needed if applications require their own customized locale definitions to
11289 operate. It is presently unknown whether many applications are dependent on this.
11290 However, the option is virtually mandatory for systems in which internationalized
11291 applications are developed.

- 11292 XSI-conformant systems support this option.
- 11293 POSIX2_PBS
- 11294 The system supports the Batch Environment option.
- 11295 POSIX2_PBS_ACCOUNTING
- 11296 The system supports the optional feature of accounting within the Batch Environment
- 11297 option. It will be required in servers that implement the optional feature of accounting.
- 11298 POSIX2_PBS_CHECKPOINT
- 11299 The system supports the optional feature of checkpoint/restart within the Batch
- 11300 Environment option.
- 11301 POSIX2_PBS_LOCATE
- 11302 The system supports the optional feature of locating batch jobs within the Batch
- 11303 Environment option.
- 11304 POSIX2_PBS_MESSAGE
- 11305 The system supports the optional feature of sending messages to batch jobs within the
- 11306 Batch Environment option.
- 11307 POSIX2_PBS_TRACK
- 11308 The system supports the optional feature of tracking batch jobs within the Batch
- 11309 Environment option.
- 11310 POSIX2_CHAR_TERM
- 11311 The system supports at least one terminal type capable of all operations described in
- 11312 IEEE Std 1003.1-2001.
- 11313 On systems with POSIX2_UPE, this option is almost always required. It was developed
- 11314 solely to allow certain specialized vendors and user applications to bypass the requirement
- 11315 for general-purpose asynchronous terminal support. For example, an application and
- 11316 system that was suitable for block-mode terminals, such as IBM 3270s, would not need this
- 11317 option.
- 11318 XSI-conformant systems support this option.

11319 D.3.3 Configurable Limits

- 11320 Very few of the limits need to be increased for profiles. No profile can cite lower values.
- 11321 {POSIX2_BC_BASE_MAX}
- 11322 {POSIX2_BC_DIM_MAX}
- 11323 {POSIX2_BC_SCALE_MAX}
- 11324 {POSIX2_BC_STRING_MAX}
- 11325 No increase is anticipated for any of these *bc* values, except for very specialized applications
- 11326 involving huge numbers.
- 11327 {POSIX2_COLL_WEIGHTS_MAX}
- 11328 Some natural languages with complex collation requirements require an increase from the
- 11329 default 2 to 4; no higher numbers are anticipated.
- 11330 {POSIX2_EXPR_NEST_MAX}
- 11331 No increase is anticipated.
- 11332 {POSIX2_LINE_MAX}
- 11333 This number is much larger than most historical applications have been able to use. At some
- 11334 future time, applications may be rewritten to take advantage of even larger values.

11335 {POSIX2_RE_DUP_MAX}
 11336 No increase is anticipated.

11337 {POSIX2_VERSION}
 11338 This is actually not a limit, but a standard version stamp. Generally, a profile should specify
 11339 the Shell and Utilities volume of IEEE Std 1003.1-2001, Chapter 2, Shell Command Language
 11340 by name in the normative references section, not this value.

11341 **D.3.4 Configuration Options (System Interfaces)**

11342 {NGROUPS_MAX}
 11343 A non-zero value indicates that the implementation supports supplementary groups.

11344 This option is needed where there is a large amount of shared use of files, but where a
 11345 certain amount of protection is needed. Many profiles³ are known to require this option; it
 11346 should only be required if needed, but it should never be prohibited.

11347 _POSIX_ADVISORY_INFO
 11348 The system provides advisory information for file management.

11349 This option allows the application to specify advisory information that can be used to
 11350 achieve better or even deterministic response time in file manager or input and output
 11351 operations.

11352 _POSIX_ASYNCHRONOUS_IO
 11353 The system provides concurrent process execution and input and output transfers.

11354 This option was created to support historical systems that did not provide the feature. It
 11355 should only be required if needed, but it should never be prohibited.

11356 _POSIX_BARRIERS
 11357 The system supports barrier synchronization.

11358 This option was created to allow efficient synchronization of multiple parallel threads in
 11359 multi-processor systems in which the operation is supported in part by the hardware
 11360 architecture.

11361 _POSIX_CHOWN_RESTRICTED
 11362 The system restricts the right to “give away” files to other users.

11363 This option should be carefully investigated before it is required. Some applications expect
 11364 that they can change the ownership of files in this way. It is provided where either security
 11365 or system account requirements cause this ability to be a problem. It is also known to be
 11366 specified in many profiles.

11367 _POSIX_CLOCK_SELECTION
 11368 The system supports the Clock Selection option.

11369 This option allows applications to request a high resolution sleep in order to suspend a
 11370 thread during a relative time interval, or until an absolute time value, using the desired
 11371 clock. It also allows the application to select the clock used in a *pthread_cond_timedwait()*
 11372 function call.

11373 _____

11374 3. There are no formally approved profiles of IEEE Std 1003.1-2001 at the time of publication; the reference here is to various
 11375 profiles generated by private bodies or governments.

11376 `_POSIX_CPUTIME`

11377 The system supports the Process CPU-Time Clocks option.

11378 This option allows applications to use a new clock that measures the execution times of
 11379 processes or threads, and the possibility to create timers based upon these clocks, for
 11380 runtime detection (and treatment) of execution time overruns.

11381 `_POSIX_FSYNC`

11382 The system supports file synchronization requests.

11383 This option was created to support historical systems that did not provide the feature.
 11384 Applications that are expecting guaranteed completion of their input and output operations
 11385 should require the `_POSIX_SYNC_IO` option. This option should never be prohibited.

11386 XSI-conformant systems support this option.

11387 `_POSIX_IPV6`

11388 The system supports facilities related to Internet Protocol Version 6 (IPv6).

11389 This option was created to allow systems to transition to IPv6.

11390 `_POSIX_JOB_CONTROL`

11391 Job control facilities are mandatory in IEEE Std 1003.1-2001.

11392 The option was created primarily to support historical systems that did not provide the
 11393 feature. Many existing profiles now require it; it should only be required if needed, but it
 11394 should never be prohibited. Most applications that use it can run when it is not present,
 11395 although with a degraded level of user convenience.

11396 `_POSIX_MAPPED_FILES`

11397 The system supports the mapping of regular files into the process address space.

11398 XSI-conformant systems support this option.

11399 Both this option and the Shared Memory Objects option provide shared access to memory
 11400 objects in the process address space. The functions defined under this option provide the
 11401 functionality of existing practice for mapping regular files. This functionality was deemed
 11402 unnecessary, if not inappropriate, for embedded systems applications and, hence, is
 11403 provided under this option. It should only be required if needed, but it should never be
 11404 prohibited.

11405 `_POSIX_MEMLOCK`

11406 The system supports the locking of the address space.

11407 This option was created to support historical systems that did not provide the feature. It
 11408 should only be required if needed, but it should never be prohibited.

11409 `_POSIX_MEMLOCK_RANGE`

11410 The system supports the locking of specific ranges of the address space.

11411 For applications that have well-defined sections that need to be locked and others that do
 11412 not, IEEE Std 1003.1-2001 supports an optional set of functions to lock or unlock a range of
 11413 process addresses. The following are two reasons for having a means to lock down a
 11414 specific range:

- 11415 1. An asynchronous event handler function that must respond to external events in a
 11416 deterministic manner such that page faults cannot be tolerated
- 11417 2. An input/output “buffer” area that is the target for direct-to-process I/O, and the
 11418 overhead of implicit locking and unlocking for each I/O call cannot be tolerated

- 11419 It should only be required if needed, but it should never be prohibited.
- 11420 `_POSIX_MEMORY_PROTECTION`
11421 The system supports memory protection.
11422 XSI-conformant systems support this option.
11423 The provision of this option typically imposes additional hardware requirements. It should
11424 never be prohibited.
- 11425 `_POSIX_PRIORITIZED_IO`
11426 The system provides prioritization for input and output operations.
11427 The use of this option may interfere with the ability of the system to optimize input and
11428 output throughput. It should only be required if needed, but it should never be prohibited.
- 11429 `_POSIX_MESSAGE_PASSING`
11430 The system supports the passing of messages between processes.
11431 This option was created to support historical systems that did not provide the feature. The
11432 functionality adds a high-performance XSI interprocess communication facility for local
11433 communication. It should only be required if needed, but it should never be prohibited.
- 11434 `_POSIX_MONOTONIC_CLOCK`
11435 The system supports the Monotonic Clock option.
11436 This option allows realtime applications to rely on a monotonically increasing clock that
11437 does not jump backwards, and whose value does not change except for the regular ticking
11438 of the clock.
- 11439 `_POSIX_PRIORITY_SCHEDULING`
11440 The system provides priority-based process scheduling.
11441 Support of this option provides predictable scheduling behavior, allowing applications to
11442 determine the order in which processes that are ready to run are granted access to a
11443 processor. It should only be required if needed, but it should never be prohibited.
- 11444 `_POSIX_REALTIME_SIGNALS`
11445 The system provides prioritized, queued signals with associated data values.
11446 This option was created to support historical systems that did not provide the features. It
11447 should only be required if needed, but it should never be prohibited.
- 11448 `_POSIX_REGEX`
11449 Support for regular expression facilities is mandatory in IEEE Std 1003.1-2001.
- 11450 `_POSIX_SAVED_IDS`
11451 Support for this feature is mandatory in IEEE Std 1003.1-2001.
11452 Certain classes of applications rely on it for proper operation, and there is no alternative
11453 short of giving the application root privileges on most implementations that did not provide
11454 `_POSIX_SAVED_IDS`.
- 11455 `_POSIX_SEMAPHORES`
11456 The system provides counting semaphores.
11457 This option was created to support historical systems that did not provide the feature. It
11458 should only be required if needed, but it should never be prohibited.
- 11459 `_POSIX_SHARED_MEMORY_OBJECTS`
11460 The system supports the mapping of shared memory objects into the process address space.

- 11461 Both this option and the Memory Mapped Files option provide shared access to memory
11462 objects in the process address space. The functions defined under this option provide the
11463 functionality of existing practice for shared memory objects. This functionality was deemed
11464 appropriate for embedded systems applications and, hence, is provided under this option. It
11465 should only be required if needed, but it should never be prohibited.
- 11466 _POSIX_SHELL
11467 Support for the *sh* utility command line interpreter is mandatory in IEEE Std 1003.1-2001.
- 11468 _POSIX_SPAWN
11469 The system supports the spawn option.
11470 This option provides applications with an efficient mechanism to spawn execution of a new
11471 process.
- 11472 _POSIX_SPINLOCKS
11473 The system supports spin locks.
11474 This option was created to support a simple and efficient synchronization mechanism for
11475 threads executing in multi-processor systems.
- 11476 _POSIX_SPORADIC_SERVER
11477 The system supports the sporadic server scheduling policy.
11478 This option provides applications with a new scheduling policy for scheduling aperiodic
11479 processes or threads in hard realtime applications.
- 11480 _POSIX_SYNCHRONIZED_IO
11481 The system supports guaranteed file synchronization.
11482 This option was created to support historical systems that did not provide the feature.
11483 Applications that are expecting guaranteed completion of their input and output operations
11484 should require this option, rather than the File Synchronization option. It should only be
11485 required if needed, but it should never be prohibited.
- 11486 _POSIX_THREADS
11487 The system supports multiple threads of control within a single process.
11488 This option was created to support historical systems that did not provide the feature.
11489 Applications written assuming a multi-threaded environment would be expected to require
11490 this option. It should only be required if needed, but it should never be prohibited.
11491 XSI-conformant systems support this option.
- 11492 _POSIX_THREAD_ATTR_STACKADDR
11493 The system supports specification of the stack address for a created thread.
11494 Applications may take advantage of support of this option for performance benefits, but
11495 dependence on this feature should be minimized. This option should never be prohibited.
11496 XSI-conformant systems support this option.
- 11497 _POSIX_THREAD_ATTR_STACKSIZE
11498 The system supports specification of the stack size for a created thread.
11499 Applications may require this option in order to ensure proper execution, but such usage
11500 limits portability and dependence on this feature should be minimized. It should only be
11501 required if needed, but it should never be prohibited.
11502 XSI-conformant systems support this option.

- 11503 `_POSIX_THREAD_PRIORITY_SCHEDULING`
11504 The system provides priority-based thread scheduling.
- 11505 Support of this option provides predictable scheduling behavior, allowing applications to
11506 determine the order in which threads that are ready to run are granted access to a processor.
11507 It should only be required if needed, but it should never be prohibited.
- 11508 `_POSIX_THREAD_PRIO_INHERIT`
11509 The system provides mutual-exclusion operations with priority inheritance.
- 11510 Support of this option provides predictable scheduling behavior, allowing applications to
11511 determine the order in which threads that are ready to run are granted access to a processor.
11512 It should only be required if needed, but it should never be prohibited.
- 11513 `_POSIX_THREAD_PRIO_PROTECT`
11514 The system supports a priority ceiling emulation protocol for mutual-exclusion operations.
- 11515 Support of this option provides predictable scheduling behavior, allowing applications to
11516 determine the order in which threads that are ready to run are granted access to a processor.
11517 It should only be required if needed, but it should never be prohibited.
- 11518 `_POSIX_THREAD_PROCESS_SHARED`
11519 The system provides shared access among multiple processes to synchronization objects.
- 11520 This option was created to support historical systems that did not provide the feature. It
11521 should only be required if needed, but it should never be prohibited.
- 11522 XSI-conformant systems support this option.
- 11523 `_POSIX_THREAD_SAFE_FUNCTIONS`
11524 The system provides thread-safe versions of all of the POSIX.1 functions.
- 11525 This option is required if the Threads option is supported. This is a separate option because
11526 thread-safe functions are useful in implementations providing other mechanisms for
11527 concurrency. It should only be required if needed, but it should never be prohibited.
- 11528 XSI-conformant systems support this option.
- 11529 `_POSIX_THREAD_SPORADIC_SERVER`
11530 The system supports the thread sporadic server scheduling policy.
- 11531 Support for this option provides applications with a new scheduling policy for scheduling
11532 aperiodic threads in hard realtime applications.
- 11533 `_POSIX_TIMEOUTS`
11534 The system provides timeouts for some blocking services.
- 11535 This option was created to provide a timeout capability to system services, thus allowing
11536 applications to include better error detection, and recovery capabilities.
- 11537 `_POSIX_TIMERS`
11538 The system provides higher resolution clocks with multiple timers per process.
- 11539 This option was created to support historical systems that did not provide the features. This
11540 option is appropriate for applications requiring higher resolution timestamps or needing to
11541 control the timing of multiple activities. It should only be required if needed, but it should
11542 never be prohibited.
- 11543 `_POSIX_TRACE`
11544 The system supports the Trace option.

- 11545 This option was created to allow applications to perform tracing.
- 11546 `_POSIX_TRACE_EVENT_FILTER`
- 11547 The system supports the Trace Event Filter option.
- 11548 This option is dependent on support of the Trace option.
- 11549 `_POSIX_TRACE_INHERIT`
- 11550 The system supports the Trace Inherit option.
- 11551 This option is dependent on support of the Trace option.
- 11552 `_POSIX_TRACE_LOG`
- 11553 The system supports the Trace Log option.
- 11554 This option is dependent on support of the Trace option.
- 11555 `_POSIX_TYPED_MEMORY_OBJECTS`
- 11556 The system supports the Typed Memory Objects option.
- 11557 This option was created to allow realtime applications to access different kinds of physical memory, and allow processes in these applications to share portions of this memory.
- 11558

11559 D.3.5 Configurable Limits

- 11560 In general, the configurable limits in the `<limits.h>` header defined in the Base Definitions
- 11561 volume of IEEE Std 1003.1-2001 have been set to minimal values; many applications or
- 11562 implementations may require larger values. No profile can cite lower values.
- 11563 `{AIO_LISTIO_MAX}`
- 11564 The current minimum is likely to be inadequate for most applications. It is expected that
- 11565 this value will be increased by profiles requiring support for list input and output
- 11566 operations.
- 11567 `{AIO_MAX}`
- 11568 The current minimum is likely to be inadequate for most applications. It is expected that
- 11569 this value will be increased by profiles requiring support for asynchronous input and
- 11570 output operations.
- 11571 `{AIO_PRIO_DELTA_MAX}`
- 11572 The functionality associated with this limit is needed only by sophisticated applications. It
- 11573 is not expected that this limit would need to be increased under a general-purpose profile.
- 11574 `{ARG_MAX}`
- 11575 The current minimum is likely to need to be increased for profiles, particularly as larger
- 11576 amounts of information are passed through the environment. Many implementations are
- 11577 believed to support larger values.
- 11578 `{CHILD_MAX}`
- 11579 The current minimum is suitable only for systems where a single user is not running
- 11580 applications in parallel. It is significantly too low for any system also requiring windows,
- 11581 and if `_POSIX_JOB_CONTROL` is specified, it should be raised.
- 11582 `{CLOCKRES_MIN}`
- 11583 It is expected that profiles will require a finer granularity clock, perhaps as fine as 1 μ s,
- 11584 represented by a value of 1 000 for this limit.
- 11585 `{DELAYTIMER_MAX}`
- 11586 It is believed that most implementations will provide larger values.

- 11587 {LINK_MAX}
11588 For most applications and usage, the current minimum is adequate. Many implementations
11589 have a much larger value, but this should not be used as a basis for raising the value unless
11590 the applications to be used require it.
- 11591 {LOGIN_NAME_MAX}
11592 This is not actually a limit, but an implementation parameter. No profile should impose a
11593 requirement on this value.
- 11594 {MAX_CANON}
11595 For most purposes, the current minimum is adequate. Unless high-speed burst serial
11596 devices are used, it should be left as is.
- 11597 {MAX_INPUT}
11598 See {MAX_CANON}.
- 11599 {MQ_OPEN_MAX}
11600 The current minimum should be adequate for most profiles.
- 11601 {MQ_PRIO_MAX}
11602 The current minimum corresponds to the required number of process scheduling priorities.
11603 Many realtime practitioners believe that the number of message priority levels ought to be
11604 the same as the number of execution scheduling priorities.
- 11605 {NAME_MAX}
11606 Many implementations now support larger values, and many applications and users
11607 assume that larger names can be used. Many existing profiles also specify a larger value.
11608 Specifying this value will reduce the number of conforming implementations, although this
11609 might not be a significant consideration over time. Values greater than 255 should not be
11610 required.
- 11611 {NGROUPS_MAX}
11612 The value selected will typically be 8 or larger.
- 11613 {OPEN_MAX}
11614 The historically common value for this has been 20. Many implementations support larger
11615 values. If applications that use larger values are anticipated, an appropriate value should be
11616 specified.
- 11617 {PAGESIZE}
11618 This is not actually a limit, but an implementation parameter. No profile should impose a
11619 requirement on this value.
- 11620 {PATH_MAX}
11621 Historically, the minimum has been either 1024 or indefinite, depending on the
11622 implementation. Few applications actually require values larger than 256, but some users
11623 may create file hierarchies that must be accessed with longer paths. This value should only
11624 be changed if there is a clear requirement.
- 11625 {PIPE_BUF}
11626 The current minimum is adequate for most applications. Historically, it has been larger. If
11627 applications that write single transactions larger than this are anticipated, it should be
11628 increased. Applications that write lines of text larger than this probably do not need it
11629 increased, as the text line is delimited by a <newline>.
- 11630 {POSIX_VERSION}
11631 This is actually not a limit, but a standard version stamp. Generally, a profile should specify
11632 IEEE Std 1003.1-2001 by a name in the normative references section, not this value.

11633	{PTHREAD_DESTRUCTOR_ITERATIONS}
11634	It is unlikely that applications will need larger values to avoid loss of memory resources.
11635	{PTHREAD_KEYS_MAX}
11636	The current value should be adequate for most profiles.
11637	{PTHREAD_STACK_MIN}
11638	This should not be treated as an actual limit, but as an implementation parameter. No
11639	profile should impose a requirement on this value.
11640	{PTHREAD_THREADS_MAX}
11641	It is believed that most implementations will provide larger values.
11642	{RTSIG_MAX}
11643	The current limit was chosen so that the set of POSIX.1 signal numbers can fit within a 32-
11644	bit field. It is recognized that most existing implementations define many more signals than
11645	are specified in POSIX.1 and, in fact, many implementations have already exceeded 32
11646	signals (including the “null signal”). Support of {_POSIX_RTSIG_MAX} additional signals
11647	may push some implementations over the single 32-bit word line, but is unlikely to push
11648	any implementations that are already over that line beyond the 64 signal line.
11649	{SEM_NSEMS_MAX}
11650	The current value should be adequate for most profiles.
11651	{SEM_VALUE_MAX}
11652	The current value should be adequate for most profiles.
11653	{SSIZE_MAX}
11654	This limit reflects fundamental hardware characteristics (the size of an integer), and should
11655	not be specified unless it is clearly required. Extreme care should be taken to assure that
11656	any value that might be specified does not unnecessarily eliminate implementations
11657	because of accidents of hardware design.
11658	{STREAM_MAX}
11659	This limit is very closely related to {OPEN_MAX}. It should never be larger than
11660	{OPEN_MAX}, but could reasonably be smaller for application areas where most files are
11661	not accessed through <i>stdio</i> . Some implementations may limit {STREAM_MAX} to 20 but
11662	allow {OPEN_MAX} to be considerably larger. Such implementations should be allowed for
11663	if the applications permit.
11664	{TIMER_MAX}
11665	The current limit should be adequate for most profiles, but it may need to be larger for
11666	applications with a large number of asynchronous operations.
11667	{TTY_NAME_MAX}
11668	This is not actually a limit, but an implementation parameter. No profile should impose a
11669	requirement on this value.
11670	{TZNAME_MAX}
11671	The minimum has been historically adequate, but if longer timezone names are anticipated
11672	(particularly such values as UTC-1), this should be increased.

11673 D.3.6 Optional Behavior

11674 In IEEE Std 1003.1-2001, there are no instances of the terms unspecified, undefined,
11675 implementation-defined, or with the verbs “may” or “need not”, that the developers of
11676 IEEE Std 1003.1-2001 anticipate or sanction as suitable for profile or test method citation. All of
11677 these are merely warnings to conforming applications to avoid certain areas that can vary from
11678 system to system, and even over time on the same system. In many cases, these terms are used
11679 explicitly to support extensions, but profiles should not anticipate and require such extensions;
11680 future versions of IEEE Std 1003.1 may do so.

11681 / *Rationale (Informative)*

11682 **Part E:**

11683 **Subprofiling Considerations**

11684 *The Open Group*

11685 *The Institute of Electrical and Electronics Engineers, Inc.*

Subprofiling Considerations (Informative)

11686

11687 This section contains further information to satisfy the requirement that the project scope enable
 11688 subprofiling of IEEE Std 1003.1-2001. The original intent was to have included a set of options
 11689 similar to the “Units of Functionality” contained in IEEE Std 1003.13-1998. However, as the
 11690 development of IEEE Std 1003.1-2001 continued, the standard developers felt it premature to fix
 11691 these in normative text. The approach instead has been to include a general requirement in
 11692 normative text regarding subprofiling and to include an informative section (here) containing a
 11693 proposed set of subprofiling options.

11694 E.1 Subprofiling Option Groups

11695 The following Option Groups⁴ are defined to support profiling. Systems claiming support to
 11696 IEEE Std 1003.1-2001 need not implement these options apart from the requirements stated in
 11697 the Base Definitions volume of IEEE Std 1003.1-2001, Section 2.1.3, POSIX Conformance. These
 11698 Option Groups allow profiles to subset the System Interfaces volume of IEEE Std 1003.1-2001 by
 11699 collecting sets of related functions.

11700 POSIX_C_LANG_JUMP: Jump Functions

11701 *longjmp()*, *setjmp()*

11702 POSIX_C_LANG_MATH: Maths Library

11703 *acos()*, *acosf()*, *acosh()*, *acoshf()*, *acoshl()*, *acosl()*, *asin()*, *asinf()*, *asinh()*, *asinhf()*, *asinhl()*,
 11704 *asinl()*, *atan()*, *atan2()*, *atan2f()*, *atan2l()*, *atanf()*, *atanh()*, *atanhf()*, *atanhl()*, *atanl()*, *cabs()*,
 11705 *cabsf()*, *cabsl()*, *cacos()*, *cacosf()*, *cacosh()*, *cacoshf()*, *cacoshl()*, *cacosl()*, *carg()*, *cargf()*, *cargl()*,
 11706 *casin()*, *casinf()*, *casinh()*, *casinhf()*, *casinhl()*, *casinl()*, *catan()*, *catanf()*, *catanh()*, *catanhf()*,
 11707 *catanhl()*, *catanl()*, *cbrt()*, *cbrtf()*, *cbrtl()*, *ccos()*, *ccosf()*, *ccosh()*, *ccoshf()*, *ccoshl()*, *ccosl()*,
 11708 *ceil()*, *ceilf()*, *ceill()*, *cexp()*, *cexpf()*, *cexpl()*, *cimag()*, *cimagf()*, *cimagl()*, *clog()*, *clogf()*, *clogl()*,
 11709 *conj()*, *conjf()*, *conjl()*, *copysign()*, *copysignf()*, *copysignl()*, *cos()*, *cosf()*, *cosh()*, *coshf()*,
 11710 *coshl()*, *cosl()*, *cpow()*, *cpowf()*, *cpowl()*, *cproj()*, *cprojf()*, *cprojl()*, *creal()*, *crealf()*, *creall()*,
 11711 *csin()*, *csinf()*, *csinh()*, *csinhf()*, *csinhl()*, *csinl()*, *csqrt()*, *csqrtf()*, *csqrtl()*, *ctan()*, *ctanf()*,
 11712 *ctanh()*, *ctanhf()*, *ctanhl()*, *ctanl()*, *erf()*, *erfc()*, *erfcf()*, *erfcl()*, *erff()*, *erfl()*, *exp()*, *exp2()*,
 11713 *exp2f()*, *exp2l()*, *expf()*, *expl()*, *expm1()*, *expm1f()*, *expm1l()*, *fabs()*, *fabsf()*, *fabsl()*, *fdim()*,
 11714 *fdimf()*, *fdiml()*, *floor()*, *floorf()*, *floorl()*, *fma()*, *fmaf()*, *fmal()*, *fmax()*, *fmaxf()*, *fmaxl()*, *fmin()*,
 11715 *fminf()*, *fminl()*, *fmod()*, *fmodf()*, *fmodl()*, *fpclassify()*, *frexp()*, *frexpf()*, *frexpl()*, *hypot()*,
 11716 *hypotf()*, *hypotl()*, *ilogb()*, *ilogbf()*, *ilogbl()*, *isfinite()*, *isgreater()*, *isgreaterequal()*, *isinf()*,
 11717 *isless()*, *islessequal()*, *islessgreater()*, *isnan()*, *isnormal()*, *isunordered()*, *ldexp()*, *ldexpf()*,
 11718 *ldexpl()*, *lgamma()*, *lgammaf()*, *lgammal()*, *llrint()*, *llrintf()*, *llrintl()*, *llround()*, *llroundf()*,
 11719 *llroundl()*, *log()*, *log10()*, *log10f()*, *log10l()*, *log1p()*, *log1pf()*, *log1pl()*, *log2()*, *log2f()*, *log2l()*,
 11720 *logb()*, *logbf()*, *logbl()*, *logf()*, *logl()*, *lrint()*, *lrintf()*, *lrintl()*, *lround()*, *lroundf()*, *lroundl()*,
 11721 *modf()*, *modff()*, *modfl()*, *nan()*, *nanf()*, *nanl()*, *nearbyint()*, *nearbyintf()*, *nearbyintl()*,
 11722 *nextafter()*, *nextafterf()*, *nextafterl()*, *nexttoward()*, *nexttowardf()*, *nexttowardl()*, *pow()*, *powf()*,
 11723 *powl()*, *remainder()*, *remainderf()*, *remainderl()*, *remquo()*, *remquof()*, *remquoil()*, *rint()*, *rintf()*,
 11724 *rintl()*, *round()*, *roundf()*, *roundl()*, *scalbln()*, *scalblnf()*, *scalblnl()*, *scalbn()*, *scalbnf()*, *scalbnl()*,
 11725 *signbit()*, *sin()*, *sinf()*, *sinh()*, *sinhf()*, *sinhl()*, *sinl()*, *sqrt()*, *sqrtf()*, *sqrtl()*, *tan()*, *tanf()*,

11726

11727 4. These are equivalent to the Units of Functionality from IEEE Std 1003.13-1998.

11728	<i>tanh()</i> , <i>tanhf()</i> , <i>tanhl()</i> , <i>tanl()</i> , <i>tgamma()</i> , <i>tgammaf()</i> , <i>tgammaL()</i> , <i>trunc()</i> , <i>truncf()</i> , <i>truncl()</i>
11729	POSIX_C_LANG_SUPPORT: General ISO C Library
11730	<i>abs()</i> , <i>asctime()</i> , <i>atof()</i> , <i>atoi()</i> , <i>atol()</i> , <i>atoll()</i> , <i>bsearch()</i> , <i>calloc()</i> , <i>ctime()</i> , <i>difftime()</i> , <i>div()</i> ,
11731	<i>feclearexcept()</i> , <i>fegetenv()</i> , <i>fegetexceptflag()</i> , <i>fegetround()</i> , <i>fehldexcept()</i> , <i>feraiseexcept()</i> ,
11732	<i>fesetenv()</i> , <i>fesetexceptflag()</i> , <i>fesetround()</i> , <i>fetestexcept()</i> , <i>feupdateenv()</i> , <i>free()</i> , <i>gmtime()</i> ,
11733	<i>imaxabs()</i> , <i>imaxdiv()</i> , <i>isalnum()</i> , <i>isalpha()</i> , <i>isblank()</i> , <i>iscntrl()</i> , <i>isdigit()</i> , <i>isgraph()</i> , <i>islower()</i> ,
11734	<i>isprint()</i> , <i>ispunct()</i> , <i>isspace()</i> , <i>isupper()</i> , <i>isxdigit()</i> , <i>labs()</i> , <i>ldiv()</i> , <i>llabs()</i> , <i>lldiv()</i> , <i>localeconv()</i> ,
11735	<i>localtime()</i> , <i>malloc()</i> , <i>memchr()</i> , <i>memcmp()</i> , <i>memcpy()</i> , <i>memmove()</i> , <i>memset()</i> , <i>mktime()</i> ,
11736	<i>qsort()</i> , <i>rand()</i> , <i>realloc()</i> , <i>setlocale()</i> , <i>snprintf()</i> , <i>sprintf()</i> , <i>srand()</i> , <i>sscanf()</i> , <i>strcat()</i> , <i>strchr()</i> ,
11737	<i>strcmp()</i> , <i>strcoll()</i> , <i>strcpy()</i> , <i>strcspn()</i> , <i>strerror()</i> , <i>strftime()</i> , <i>strlen()</i> , <i>strncat()</i> , <i>strncpy()</i> ,
11738	<i>strncpy()</i> , <i>strpbrk()</i> , <i>strrchr()</i> , <i>strspn()</i> , <i>strstr()</i> , <i>strtod()</i> , <i>strtof()</i> , <i>strtoimax()</i> , <i>strtok()</i> , <i>strtol()</i> ,
11739	<i>strtold()</i> , <i>strtoll()</i> , <i>strtoul()</i> , <i>strtoull()</i> , <i>strtoumax()</i> , <i>strxfrm()</i> , <i>time()</i> , <i>tolower()</i> , <i>toupper()</i> ,
11740	<i>tzname()</i> , <i>tzset()</i> , <i>va_arg()</i> , <i>va_copy()</i> , <i>va_end()</i> , <i>va_start()</i> , <i>vsprintf()</i> , <i>vsscanf()</i>
11741	POSIX_C_LANG_SUPPORT_R: Thread-Safe General ISO C Library
11742	<i>asctime_r()</i> , <i>ctime_r()</i> , <i>gmtime_r()</i> , <i>localtime_r()</i> , <i>rand_r()</i> , <i>strerror_r()</i> , <i>strtok_r()</i>
11743	POSIX_C_LANG_WIDE_CHAR: Wide-Character ISO C Library
11744	<i>btowc()</i> , <i>iswalnum()</i> , <i>iswalpha()</i> , <i>iswblank()</i> , <i>iswcntrl()</i> , <i>iswctype()</i> , <i>iswdigit()</i> , <i>iswgraph()</i> ,
11745	<i>iswlower()</i> , <i>iswprint()</i> , <i>iswpunct()</i> , <i>iswspace()</i> , <i>iswupper()</i> , <i>iswxdigit()</i> , <i>mblen()</i> , <i>mbrlen()</i> ,
11746	<i>mbrtowc()</i> , <i>mbsinit()</i> , <i>mbsrtowcs()</i> , <i>mbstowcs()</i> , <i>mbtowc()</i> , <i>swprintf()</i> , <i>swscanf()</i> , <i>towctrans()</i> ,
11747	<i>towlower()</i> , <i>towupper()</i> , <i>vswprintf()</i> , <i>vswscanf()</i> , <i>wcrtomb()</i> , <i>wcscat()</i> , <i>wcchr()</i> , <i>wcscmp()</i> ,
11748	<i>wcscoll()</i> , <i>wcscpy()</i> , <i>wcscspn()</i> , <i>wcsftime()</i> , <i>wcslen()</i> , <i>wcsncat()</i> , <i>wcsncmp()</i> , <i>wcsncpy()</i> ,
11749	<i>wcspbrk()</i> , <i>wcsrchr()</i> , <i>wcrtombs()</i> , <i>wcsspn()</i> , <i>wcsstr()</i> , <i>wcstod()</i> , <i>wcstof()</i> , <i>wcstoimax()</i> ,
11750	<i>wcstok()</i> , <i>wcstol()</i> , <i>wcstold()</i> , <i>wcstoll()</i> , <i>wcstombs()</i> , <i>wcstoul()</i> , <i>wcstoull()</i> , <i>wcstoumax()</i> ,
11751	<i>wcsxfrm()</i> , <i>wctob()</i> , <i>wctomb()</i> , <i>wctrans()</i> , <i>wctype()</i> , <i>wmemchr()</i> , <i>wmemcmp()</i> , <i>wmemcpy()</i> ,
11752	<i>wmemmove()</i> , <i>wmemset()</i>
11753	POSIX_C_LIB_EXT: General C Library Extension
11754	<i>fnmatch()</i> , <i>getopt()</i> , <i>optarg</i> , <i>opterr</i> , <i>optind</i> , <i>optopt</i>
11755	POSIX_DEVICE_IO: Device Input and Output
11756	<i>FD_CLR()</i> , <i>FD_ISSET()</i> , <i>FD_SET()</i> , <i>FD_ZERO()</i> , <i>clearerr()</i> , <i>close()</i> , <i>fclose()</i> , <i>fdopen()</i> , <i>feof()</i> ,
11757	<i>ferror()</i> , <i>fflush()</i> , <i>fgetc()</i> , <i>fgets()</i> , <i>fileno()</i> , <i>fopen()</i> , <i>fprintf()</i> , <i>fputc()</i> , <i>fputs()</i> , <i>fread()</i> , <i>freopen()</i> ,
11758	<i>fscanf()</i> , <i>fwrite()</i> , <i>getc()</i> , <i>getchar()</i> , <i>gets()</i> , <i>open()</i> , <i>perror()</i> , <i>printf()</i> , <i>pselect()</i> , <i>putc()</i> , <i>putchar()</i> ,
11759	<i>puts()</i> , <i>read()</i> , <i>scanf()</i> , <i>select()</i> , <i>setbuf()</i> , <i>setvbuf()</i> , <i>stderr</i> , <i>stdin</i> , <i>stdout</i> , <i>ungetc()</i> , <i>vfprintf()</i> ,
11760	<i>vfscanf()</i> , <i>vprintf()</i> , <i>vscanf()</i> , <i>write()</i>
11761	POSIX_DEVICE_SPECIFIC: General Terminal
11762	<i>cfgetispeed()</i> , <i>cfgetospeed()</i> , <i>cfsetispeed()</i> , <i>cfsetospeed()</i> , <i>ctermid()</i> , <i>isatty()</i> , <i>tcdrain()</i> , <i>tclflow()</i> ,
11763	<i>tcflush()</i> , <i>tcgetattr()</i> , <i>tcsendbreak()</i> , <i>tcsetattr()</i> , <i>ttyname()</i>
11764	POSIX_DEVICE_SPECIFIC_R: Thread-Safe General Terminal
11765	<i>ttyname_r()</i>
11766	POSIX_FD_MGMT: File Descriptor Management
11767	<i>dup()</i> , <i>dup2()</i> , <i>fcntl()</i> , <i>fgetpos()</i> , <i>fseek()</i> , <i>fseeko()</i> , <i>fsetpos()</i> , <i>ftell()</i> , <i>ftello()</i> , <i>ftruncate()</i> , <i>lseek()</i> ,
11768	<i>rewind()</i>
11769	POSIX_FIFO: FIFO
11770	<i>mkfifo()</i>
11771	POSIX_FILE_ATTRIBUTES: File Attributes
11772	<i>chmod()</i> , <i>chown()</i> , <i>fchmod()</i> , <i>fchown()</i> , <i>umask()</i>
11773	POSIX_FILE_LOCKING: Thread-Safe Stdio Locking
11774	<i>flockfile()</i> , <i>ftrylockfile()</i> , <i>funlockfile()</i> , <i>getc_unlocked()</i> , <i>getchar_unlocked()</i> , <i>putc_unlocked()</i> ,

11775	<i>putchar_unlocked()</i>
11776	POSIX_FILE_SYSTEM: File System
11777	<i>access()</i> , <i>chdir()</i> , <i>closedir()</i> , <i>creat()</i> , <i>fpathconf()</i> , <i>fstat()</i> , <i>getcwd()</i> , <i>link()</i> , <i>mkdir()</i> , <i>opendir()</i> ,
11778	<i>pathconf()</i> , <i>readdir()</i> , <i>remove()</i> , <i>rename()</i> , <i>rewinddir()</i> , <i>rmdir()</i> , <i>stat()</i> , <i>tmpfile()</i> , <i>tmpnam()</i> ,
11779	<i>unlink()</i> , <i>utime()</i>
11780	POSIX_FILE_SYSTEM_EXT: File System Extensions
11781	<i>glob()</i> , <i>globfree()</i>
11782	POSIX_FILE_SYSTEM_R: Thread-Safe File System
11783	<i>readdir_r()</i>
11784	POSIX_JOB_CONTROL: Job Control
11785	<i>setpgid()</i> , <i>tcgetpgrp()</i> , <i>tcsetpgrp()</i>
11786	POSIX_MULTI_PROCESS: Multiple Processes
11787	<i>_Exit()</i> , <i>_exit()</i> , <i>assert()</i> , <i>atexit()</i> , <i>clock()</i> , <i>execl()</i> , <i>execle()</i> , <i>execlp()</i> , <i>execv()</i> , <i>execve()</i> , <i>execvp()</i> ,
11788	<i>exit()</i> , <i>fork()</i> , <i>getpgrp()</i> , <i>getpid()</i> , <i>getppid()</i> , <i>setsid()</i> , <i>sleep()</i> , <i>times()</i> , <i>wait()</i> , <i>waitpid()</i>
11789	POSIX_NETWORKING: Networking
11790	<i>accept()</i> , <i>bind()</i> , <i>connect()</i> , <i>endhostent()</i> , <i>endnetent()</i> , <i>endprotoent()</i> , <i>endservent()</i> ,
11791	<i>freeaddrinfo()</i> , <i>gai_strerror()</i> , <i>getaddrinfo()</i> , <i>gethostbyaddr()</i> , <i>gethostbyname()</i> , <i>gethostent()</i> ,
11792	<i>gethostname()</i> , <i>getnameinfo()</i> , <i>getnetbyaddr()</i> , <i>getnetbyname()</i> , <i>getnetent()</i> , <i>getpeername()</i> ,
11793	<i>getprotobyname()</i> , <i>getprotobynumber()</i> , <i>getprotoent()</i> , <i>getservbyname()</i> , <i>getservbyport()</i> ,
11794	<i>getservent()</i> , <i>getsockname()</i> , <i>getsockopt()</i> , <i>h_errno</i> , <i>htonl()</i> , <i>htons()</i> , <i>if_freenameindex()</i> ,
11795	<i>if_indextoname()</i> , <i>if_nameindex()</i> , <i>if_nametoindex()</i> , <i>inet_addr()</i> , <i>inet_ntoa()</i> , <i>inet_ntop()</i> ,
11796	<i>inet_pton()</i> , <i>listen()</i> , <i>ntohl()</i> , <i>ntohs()</i> , <i>recv()</i> , <i>recvfrom()</i> , <i>recvmsg()</i> , <i>send()</i> , <i>sendmsg()</i> , <i>sendto()</i> ,
11797	<i>sethostent()</i> , <i>setnetent()</i> , <i>setprotoent()</i> , <i>setservent()</i> , <i>setsockopt()</i> , <i>shutdown()</i> , <i>socket()</i> ,
11798	<i>socketatmark()</i> , <i>socketpair()</i>
11799	POSIX_PIPE: Pipe
11800	<i>pipe()</i>
11801	POSIX_REGEX: Regular Expressions
11802	<i>regcomp()</i> , <i>regerror()</i> , <i>regexexec()</i> , <i>regfree()</i>
11803	POSIX_SHELL_FUNC: Shell and Utilities
11804	<i>pclose()</i> , <i>popen()</i> , <i>system()</i> , <i>wordexp()</i> , <i>wordfree()</i>
11805	POSIX_SIGNALS: Signal
11806	<i>abort()</i> , <i>alarm()</i> , <i>kill()</i> , <i>pause()</i> , <i>raise()</i> , <i>sigaction()</i> , <i>sigaddset()</i> , <i>sigdelset()</i> , <i>sigemptyset()</i> ,
11807	<i>sigfillset()</i> , <i>sigismember()</i> , <i>signal()</i> , <i>sigpending()</i> , <i>sigprocmask()</i> , <i>sigsuspend()</i> , <i>sigwait()</i>
11808	POSIX_SIGNAL_JUMP: Signal Jump Functions
11809	<i>siglongjmp()</i> , <i>sigsetjmp()</i>
11810	POSIX_SINGLE_PROCESS: Single Process
11811	<i>confstr()</i> , <i>environ</i> , <i>errno</i> , <i>getenv()</i> , <i>setenv()</i> , <i>sysconf()</i> , <i>uname()</i> , <i>unsetenv()</i>
11812	POSIX_SYMBOLIC_LINKS: Symbolic Links
11813	<i>lstat()</i> , <i>readlink()</i> , <i>symlink()</i>
11814	POSIX_SYSTEM_DATABASE: System Database
11815	<i>getgrgid()</i> , <i>getgrnam()</i> , <i>getpwnam()</i> , <i>getpwuid()</i>
11816	POSIX_SYSTEM_DATABASE_R: Thread-Safe System Database
11817	<i>getgrgid_r()</i> , <i>getgrnam_r()</i> , <i>getpwnam_r()</i> , <i>getpwuid_r()</i>

11818	POSIX_USER_GROUPS: User and Group
11819	<i>getegid()</i> , <i>geteuid()</i> , <i>getgid()</i> , <i>getgroups()</i> , <i>getlogin()</i> , <i>getuid()</i> , <i>setegid()</i> , <i>seteuid()</i> , <i>setgid()</i> ,
11820	<i>setuid()</i>
11821	POSIX_USER_GROUPS_R: Thread-Safe User and Group
11822	<i>getlogin_r()</i>
11823	POSIX_WIDE_CHAR_DEVICE_IO: Device Input and Output
11824	<i>fgetwc()</i> , <i>fgetws()</i> , <i>fputwc()</i> , <i>fputws()</i> , <i>fwide()</i> , <i>fwprintf()</i> , <i>fwscanf()</i> , <i>getwc()</i> , <i>getwchar()</i> ,
11825	<i>putwc()</i> , <i>putwchar()</i> , <i>ungetwc()</i> , <i>vfwprintf()</i> , <i>vfwscanf()</i> , <i>vwprintf()</i> , <i>vwscanf()</i> , <i>wprintf()</i> ,
11826	<i>wscanf()</i>
11827	XSI_C_LANG_SUPPORT: XSI General C Library
11828	<i>_tolower()</i> , <i>_toupper()</i> , <i>a64l()</i> , <i>daylight()</i> , <i>drand48()</i> , <i>erand48()</i> , <i>ffs()</i> , <i>getcontext()</i> , <i>getdate()</i> ,
11829	<i>getsubopt()</i> , <i>hcreate()</i> , <i>hdestroy()</i> , <i>hsearch()</i> , <i>iconv()</i> , <i>iconv_close()</i> , <i>iconv_open()</i> , <i>initstate()</i> ,
11830	<i>insque()</i> , <i>isascii()</i> , <i>jrand48()</i> , <i>l64a()</i> , <i>lcong48()</i> , <i>lfind()</i> , <i>lrand48()</i> , <i>lsearch()</i> , <i>makecontext()</i> ,
11831	<i>memccpy()</i> , <i>mrand48()</i> , <i>nrand48()</i> , <i>random()</i> , <i>remque()</i> , <i>seed48()</i> , <i>setcontext()</i> , <i>setstate()</i> ,
11832	<i>signgam</i> , <i>srand48()</i> , <i>srandom()</i> , <i>strcasecmp()</i> , <i>strdup()</i> , <i>strfmon()</i> , <i>strncasecmp()</i> , <i>strptime()</i> ,
11833	<i>swab()</i> , <i>swapcontext()</i> , <i>tdelete()</i> , <i>tfind()</i> , <i>timezone()</i> , <i>toascii()</i> , <i>tsearch()</i> , <i>twalk()</i>
11834	XSI_DBM: XSI Database Management
11835	<i>dbm_clearerr()</i> , <i>dbm_close()</i> , <i>dbm_delete()</i> , <i>dbm_error()</i> , <i>dbm_fetch()</i> , <i>dbm_firstkey()</i> ,
11836	<i>dbm_nextkey()</i> , <i>dbm_open()</i> , <i>dbm_store()</i>
11837	XSI_DEVICE_IO: XSI Device Input and Output
11838	<i>fntmsg()</i> , <i>poll()</i> , <i>pread()</i> , <i>pwrite()</i> , <i>readv()</i> , <i>writenv()</i>
11839	XSI_DEVICE_SPECIFIC: XSI General Terminal
11840	<i>grantpt()</i> , <i>posix_openpt()</i> , <i>ptsname()</i> , <i>unlockpt()</i>
11841	XSI_DYNAMIC_LINKING: XSI Dynamic Linking
11842	<i>dlclose()</i> , <i>dlerror()</i> , <i>dlopen()</i> , <i>dlsym()</i>
11843	XSI_FD_MGMT: XSI File Descriptor Management
11844	<i>truncate()</i>
11845	XSI_FILE_SYSTEM: XSI File System
11846	<i>basename()</i> , <i>dirname()</i> , <i>fchdir()</i> , <i>fstatvfs()</i> , <i>ftw()</i> , <i>lchown()</i> , <i>lockf()</i> , <i>mknod()</i> , <i>mkstemp()</i> , <i>nftw()</i> ,
11847	<i>realpath()</i> , <i>seekdir()</i> , <i>statvfs()</i> , <i>sync()</i> , <i>telldir()</i> , <i>tempnam()</i>
11848	XSI_I18N: XSI Internationalization
11849	<i>catclose()</i> , <i>catgets()</i> , <i>catopen()</i> , <i>nl_langinfo()</i>
11850	XSI_IPC: XSI Interprocess Communication
11851	<i>flok()</i> , <i>msgctl()</i> , <i>msgget()</i> , <i>msgrcv()</i> , <i>msgsnd()</i> , <i>semctl()</i> , <i>semget()</i> , <i>semop()</i> , <i>shmat()</i> , <i>shmctl()</i> ,
11852	<i>shmdt()</i> , <i>shmget()</i>
11853	XSI_JOB_CONTROL: XSI Job Control
11854	<i>tcgetsid()</i>
11855	XSI_JUMP: XSI Jump Functions
11856	<i>_longjmp()</i> , <i>_setjmp()</i>
11857	XSI_MATH: XSI Maths Library
11858	<i>j0()</i> , <i>j1()</i> , <i>jn()</i> , <i>scalb()</i> , <i>y0()</i> , <i>y1()</i> , <i>yn()</i>
11859	XSI_MULTI_PROCESS: XSI Multiple Process
11860	<i>getpgid()</i> , <i>getpriority()</i> , <i>getrlimit()</i> , <i>getrusage()</i> , <i>getsid()</i> , <i>nice()</i> , <i>setpgrp()</i> , <i>setpriority()</i> ,
11861	<i>setrlimit()</i> , <i>ulimit()</i> , <i>usleep()</i> , <i>vfork()</i> , <i>waitid()</i>

11862	XSI_SIGNALS: XSI Signal
11863	<i>bsd_signal()</i> , <i>killpg()</i> , <i>sigaltstack()</i> , <i>sighold()</i> , <i>sigignore()</i> , <i>siginterrupt()</i> , <i>sigpause()</i> , <i>sigrelse()</i> ,
11864	<i>sigset()</i> , <i>ualarm()</i>
11865	XSI_SINGLE_PROCESS: XSI Single Process
11866	<i>gethostid()</i> , <i>gettimeofday()</i> , <i>putenv()</i>
11867	XSI_SYSTEM_DATABASE: XSI System Database
11868	<i>endpwent()</i> , <i>getpwent()</i> , <i>setpwent()</i>
11869	XSI_SYSTEM_LOGGING: XSI System Logging
11870	<i>closelog()</i> , <i>openlog()</i> , <i>setlogmask()</i> , <i>syslog()</i>
11871	XSI_THREAD_MUTEX_EXT: XSI Thread Mutex Extensions
11872	<i>pthread_mutexattr_gettype()</i> , <i>pthread_mutexattr_settype()</i>
11873	XSI_THREADS_EXT: XSI Threads Extensions
11874	<i>pthread_attr_getguardsize()</i> , <i>pthread_attr_getstack()</i> , <i>pthread_attr_setguardsize()</i> ,
11875	<i>pthread_attr_setstack()</i> , <i>pthread_getconcurrency()</i> , <i>pthread_setconcurrency()</i>
11876	XSI_TIMERS: XSI Timers
11877	<i>getitimer()</i> , <i>setitimer()</i>
11878	XSI_USER_GROUPS: XSI User and Group
11879	<i>endgrent()</i> , <i>endutxent()</i> , <i>getgrent()</i> , <i>getutxent()</i> , <i>getutxid()</i> , <i>getutxline()</i> , <i>pututxline()</i> ,
11880	<i>setgrent()</i> , <i>setregid()</i> , <i>setreuid()</i> , <i>setutxent()</i>
11881	XSI_WIDE_CHAR: XSI Wide-Character Library
11882	<i>wcswidth()</i> , <i>wcwidth()</i>

Index

/dev/tty.....	21
/etc/passwd	34
<pthread.h>	162
_asm_builtin_atoi().....	86
_exit()	102, 120, 271
_Exit()	271
_longjmp()	274
_POSIX_ADVISORY_INFO	280
_POSIX_ASYNCHRONOUS_IO	280
_POSIX_BARRIERS.....	280
_POSIX_CHOWN_RESTRICTED	4, 280
_POSIX_CLOCK_SELECTION	280
_POSIX_CPUTIME.....	281
_POSIX_C_SOURCE.....	87, 91
_POSIX_FSYNC	281
_POSIX_IPV6.....	281
_POSIX_JOB_CONTROL.....	5, 281, 285
_POSIX_MAPPED_FILES.....	281
_POSIX_MEMLOCK.....	281
_POSIX_MEMLOCK_RANGE	281
_POSIX_MEMORY_PROTECTION	282
_POSIX_MESSAGE_PASSING.....	282
_POSIX_MONOTONIC_CLOCK.....	282
_POSIX_NO_TRUNC	5
_POSIX_PRIORITIZED_IO	282
_POSIX_PRIORITY_SCHEDULING.....	282
_POSIX_REALTIME_SIGNALS	282
_POSIX_REGEX	282
_POSIX_RTSIG_MAX	96, 287
_POSIX_SAVED_IDS.....	5, 282
_POSIX_SEMAPHORES	282
_POSIX_SHARED_MEMORY_OBJECTS	282
_POSIX_SHELL.....	283
_POSIX_SOURCE.....	87
_POSIX_SPAWN	283
_POSIX_SPINLOCKS	283
_POSIX_SPORADIC_SERVER	283
_POSIX_SS_REPL_MAX.....	135
_POSIX_SYNCHRONIZED_IO	283
_POSIX_SYNC_IO	281
_POSIX_THREADS.....	283
_POSIX_THREAD_ATTR_STACKADDR.....	283
_POSIX_THREAD_ATTR_STACKSIZE.....	283
_POSIX_THREAD_PRIORITY_SCHEDULING	284
.....	284
_POSIX_THREAD_PRIO_INHERIT	284
.....	284
_POSIX_THREAD_PRIO_PROTECT	284
.....	284
_POSIX_THREAD_PROCESS_SHARED.....	284
.....	284
_POSIX_THREAD_SAFE_FUNCTIONS.....	284
.....	284
_POSIX_THREAD_SPORADIC_SERVER	284
.....	284
_POSIX_TIMEOUTS	284
.....	284
_POSIX_TIMERS	284
.....	284
_POSIX_TRACE.....	284
_POSIX_TRACE_EVENT_FILTER.....	285
_POSIX_TRACE_INHERIT.....	285
_POSIX_TRACE_LOG.....	285
_POSIX_TYPED_MEMORY_OBJECTS	285
_POSIX_TZNAME_MAX	58
_POSIX_VDISABLE	5
_SC_PAGESIZE	119, 121
_setjmp()	274
_XOPEN_SOURCE	87
__errno().....	94
abort()	271
access()	33, 272
active trace stream.....	204
adb, rationale for omission.....	260
address families	178
addressing.....	178
advisory information.....	106
aio_cancel().....	114-115
aio_fsync()	99, 113
AIO_LISTIO_MAX.....	285
AIO_MAX	285
AIO_PRIO_DELTA_MAX.....	285
aio_read().....	115
aio_suspend()	114, 139
aio_write()	115
alarm()	104, 138, 271
alias	229, 274
alias substitution.....	233
AND lists.....	248
application instrumentation.....	192
appropriate privilege.....	13, 25
ar	277-278
arbitrary file size	228
ARG_MAX	23, 222, 285
arithmetic expansion	239
as, rationale for omission.....	260
asa	275, 277-278
ASCII.....	24
async-cancel safety.....	176

- async-signal-safe.....102
- asynchronous error179
- asynchronous I/O.....113, 272
- asynchronous lists.....247
- at275
- atexit()175
- atoi().....86
- awk.....275, 277
- background19-21, 68
- backslash231
- banner, rationale for omission260
- barriers.....155
- basename274-275
- basic regular expression.....60
- batch275
 - general concepts257
- batch environment254
 - option definitions255
- batch environment utilities
 - common behavior260
- batch services259
- batch systems
 - historical implementations.....254
 - history254
- bc.....274-275, 279
- BC_BASE_MAX223
- BC_DIM_MAX223
- BC_SCALE_MAX223
- bg.....229, 275
- bounded response273
- bracket expression
 - grammar.....65
- BRE
 - expression anchoring.....62
 - grammar lexical conventions.....64
 - matching a collating element.....60
 - matching a single character.....60
 - matching multiple characters62
 - ordinary character60
 - periods60
 - precedence62
 - special character60
- BSD.....16, 18, 21, 24, 67-69, 95-98, 101-102, 205
- built-in utilities.....229
- C Shell.....18-20
- C-language extensions269, 274
- c99.....277
- calendar, rationale for omission260
- cancel, rationale for omission260
- cancelation cleanup handler.....174-175
- cancelation cleanup stack174
- canonical mode input processing69
- case248
- case folding.....34-35
- cat.....228
- catclose()276
- catgets().....276
- catopen()276
- cd230, 275
- CEO61
- change history.....79, 219
- character14
- character encoding42
- character set42
- character set description file.....42
- character set, portable filename.....24
- character, rationale14
- CHARCLASS_NAME_MAX47
- CHAR_MAX.....49
- chgrp228, 275-276
- CHILD_MAX.....5, 206, 223, 247, 285
- chmod228, 272, 275-276
- chmod()26
- chown228, 272, 275-276
- chown()26
- chroot().....25
- chroot, rationale for omission260
- cksum228, 275-276
- clock.....135
- clock tick.....14
- clock tick, rationale.....14
- CLOCKRES_MIN285
- clocks135
- clock_getcpuclockid().....141, 143
- clock_nanosleep().....139-140
- CLOCK_PROCESS_CPUTIME_ID.....141, 143
- CLOCK_REALTIME135-138, 140
- CLOCK_THREAD_CPUTIME_ID141, 143
- close()120, 272
- closedir()272
- closelog().....276
- cmp.....228, 275
- codes8, 85, 221
- col, rationale for omission260
- collating element order.....61
- COLL_WEIGHTS_MAX223
- column position15
- COLUMNS.....57
- comm.....275
- command14, 230, 274
- command execution.....245
- command language269, 274

Index

command search.....	245	directory device	65
command substitution	238	directory entry.....	16
compilation environment	87	directory files.....	65
complex data manipulation.....	269, 275	directory protection	33
compound commands.....	248	directory structure.....	65
concurrent execution	33	directory, root.....	25
conditional construct		dirname.....	274-275
case	248	dis, rationale for omission	261
if	249	display.....	16
configurable limits	279, 285	dot	16
configuration interrogation.....	268, 271	dot-dot	16, 24, 38
configuration options	277	double-quotes.....	231
shell and utilities	278	du	228, 275-276
system interfaces	280	dup().....	120, 272
conformance	5, 9, 12-13, 35, 79, 204, 219	dup2().....	120, 272
conformance document.....	5	EBUSY.....	93, 162
conformance document, rationale	5	ECANCELED	92
conforming application.....	12, 95, 226-227	echo.....	274
conforming application, strictly.....	9, 12, 102	ECHOE	71
confstr	272	ECHOK	71
connection indication queue	179	ECHONL.....	71
control mode.....	71	ed.....	275-276
controlling terminal.....	15, 68, 271	EDOM	93
core	34	EFAULT	91-92
covert channel	35	effective user ID	33
cp	228, 275	EFTYPE	92
cpio, rationale for omission	260	EILSEQ.....	93
cpp, rationale for omission.....	260	EINPROGRESS	114
creat()	120, 228	EINTR.....	92, 94, 104
crontab	275	EINVAL	92
CSIZE	71	ELOOP	92
csplit	275	emacs, rationale for omission	261
ctags.....	277	ENAMETOOLONG	93
ctime()	273	endgrent()	32
cu, rationale for omission	260	endpwent()	32
cut	275	ENOMEM.....	93
data access	268, 272	ENOSYS	93, 117
data type.....	204	ENOTSUP	93
date	276	ENOTTY	66, 92-93
dc, rationale for omission	261	env	274, 276
dd.....	228, 275	environment access.....	268, 272
definitions	85, 220	environment variable.....	55
DELAYTIMER_MAX	285	definition	55
determinism.....	268	EOVERFLOW.....	93
device number.....	15	EPERM.....	171
device, logical.....	22	EPIPE.....	93
df	228, 275-276	Epoch.....	17, 39, 135
diff.....	275	ERANGE.....	93
dircmp, rationale for omission	261	ERASE.....	69
direct I/O	16	ERE	63
directory	16	alternation.....	64

- bracket expression.....64
- expression anchoring.....64
- grammar.....65
- grammar lexical conventions.....64
- matching a collating element.....63
- matching a single character.....63
- matching multiple characters.....64
- ordinary character.....63
- periods.....63
- precedence.....64
- special character.....63
- EROFS.....93
- errno.....91
 - per-thread.....94
- error conditions.....41
- error handling.....252
- error numbers.....91, 95
- errors.....243
- escape character.....231
- ex.....274-276
- exec.....198
- exec family.....19, 118, 175, 224, 229, 250, 271
- Execution Time Monitoring.....141
- execution time, measurement.....36
- exit status.....243
- exit().....102, 271
- EXIT_FAILURE.....273
- EXIT_SUCCESS.....273
- expand.....275
- expr.....274-275
- EXPR_NEST_MAX.....223
- extended regular expression.....63
- extended security controls.....33
- false.....230, 274
- fc.....229, 274
- fchmod.....272
- fclose.....272
- fcntl().....66, 92, 120, 123, 205, 272
- fcntl() locks.....177
- fdopen().....120
- FD_CLOEXEC.....123
- feature test macro.....87-88, 205
- fg.....229, 275
- fgetc.....272
- fgetpos().....205
- field splitting.....241
- FIFO.....17, 24, 129
- FIFO special file.....17
- file.....17
- file access permissions.....33
- file classes.....17
- file descriptors.....104
- file format notation.....41
- file hierarchy.....34
- file hierarchy manipulation.....269, 275
- file permissions.....33, 68
- file removal.....221
- file size, arbitrary.....228
- file system.....17
- file system, mounted.....22
- file system, root.....25
- file times update.....35
- file, passwd.....24
- filename.....17, 34
- fileno().....23
- filtering of trace event types.....201
- filtering trace event types.....201
- find.....275-276
- FIPS requirements.....4
- flockfile().....94
- fnmatch().....277
- fold.....275
- fopen().....18, 228, 272
- for loop.....248
- foreground.....19-21, 67-68
- fork().....19, 67, 109, 118, 120, 198, 205, 207, 271
- fort77.....277
- fpathconf().....271-2
- fputc().....272
- fread().....272
- free().....103, 174
- fseek().....272
- fseeko().....205
- fsetpos().....205, 272
- fstat().....271-272
- fsync().....113
- ftello().....205
- ftruncate().....119-120, 122, 272
- ftw().....228
- function definition command.....249
- functions
 - implementation of.....86
 - use of.....86
- fwrite().....272
- genccat.....276
- general terminal interface.....66
- getc().....165, 272
- getch.....272
- getconf.....222, 271, 276-277
- getegid().....271
- geteuid().....271
- getgid().....271

Index

- getgrent()32
- getgrgid().....32, 166, 271
- getgrnam()32, 35, 166, 271
- getgroups()26
- getlogin().....271
- getopt()74, 276-277
- getopts230, 277
- getpgrp()20
- getpid()104, 271
- getppid()271
- getpriority().....130
- getpwent()32
- getpwnam()32, 35, 166, 271
- getpwuid()32, 166, 271
- getrlimit().....229
- getrusage().....143
- getty68
- getuid()104, 206, 271
- gid_t.....32
- glob().....277
- globfree().....277
- gmtime()40
- grammar conventions.....224
- grep275-276
- group database.....17
- group database access32
- group file18
- grouping commands.....248
- head275
- headers.....75
- here-document243
- historical implementations.....18
- HOME.....5
- host byte order36
- hosted implementation18
- ICANON69, 71
- iconv()276
- iconv_close()276
- iconv_open()276
- id276
- if249
- implementation.....18
- implementation, historical.....18
- implementation, hosted18
- implementation, native23
- implementation, specific18
- implementation-defined5-7
- implementation-defined, rationale5
- incomplete pathname18
- input file descriptor
 - duplication.....243
- input mode.....70
- input processing69
 - canonical mode69
 - non-canonical mode.....69
- inter-user communication.....270, 276
- interactive facilities.....269, 274
- interface characteristics.....67
- interfaces178
- internationalization variable.....56
- interprocess communication.....105
- invalid
 - use in RE.....60
- ioctl().....66, 93
- IPC105
- ISO C standard.....12, 14, 40, 66, 86-88
 -90-91, 93, 95, 102, 205
- ISO/IEC 646: 1991 standard24
- ISTRIP70
- job control.....18-21, 23-24, 67-68, 95, 102, 271
- job control, implementing applications.....21
- job control, implementing shells19
- job control, implementing systems.....21
- jobs229, 275
- join275
- kernel.....22
- kernel entity.....168
- kill.....230, 275
- kill().....95-97, 99-100, 102, 205
- last close120
- lchown()26
- LC_COLLATE47
- LC_CTYPE46
- LC_MESSAGES.....51
- LC_MONETARY49
- LC_NUMERIC50, 74
- LC_TIME50
- ld, rationale for omission261
- legacy.....6
- legacy, rationale6
- lex277-278
- library routine.....22
- limits.....222
- line, rationale for omission261
- LINES57
- LINE_MAX.....23, 31, 223
- link272
- link()16, 26
- LINK_MAX223, 286
- lint, rationale for omission.....261
- lio_listio().....99, 114
- lists.....247

- In228, 275
- local mode71
- locale45, 276
 - definition example52
 - definition grammar52
 - grammar52
 - lexical conventions52
- locale configuration269, 276
- locale definition45
- localedef276-277
- localtime()40, 165
- logger276
- logical device22
- login, rationale for omission261
- LOGIN_NAME_MAX286
- LOGNAME5, 57
- logname276
- longjmp()92, 103, 172, 174, 274
- LONG_MAX72
- LONG_MIN72
- lorder, rationale for omission261
- lp276
- lpstat, rationale for omission261
- ls228, 275
- lseek()113-114, 120, 164, 205, 272
- lstat()26, 228, 271-2
- lutime()26
- macro7
- mail, rationale for omission261
- mailx274, 276
- make277-278
- malloc()103, 123, 125, 153, 165-166, 174, 176
- man274
- map22
- mapped22
- margin code notation8
- margin codes8, 85, 221
- mathematical functions
 - error conditions41
 - NaN arguments41
- MAX_CANON69, 223, 286
- MAX_INPUT223, 286
- may5
- may, rationale5
- MCL_FUTURE117
- MCL_INHERIT118
- memory locking116
- memory management116, 272
- memory management unit117
- memory object22
- memory synchronization36
- memory-resident22
- mesg276
- message passing107, 272-273
- message queue107
- mkdir272, 275
- mkfifo275
- mkfifo()18
- mknod()18
- mknod, rationale for omission261
- mktime()40
- mlockall()117
- mmap()119-121, 123
- MMU117
- modem disconnect70
- monotonic clock139
- more274, 276
- mount point22
- mount()22
- mounted file system22
- mprotect()120
- mq_open()108
- MQ_OPEN_MAX286
- MQ_PRIO_MAX286
- mq_receive()109
- mq_send()109
- mq_timedreceive()140
- mq_timedsend()140
- msg*()105
- msgctl()105
- msgget()105
- msgrcv()105
- msgsnd()105
- msync()120
- multi-byte character17, 69-70
- multiple tasks269, 275
- munmap()119-121, 123, 126
- mutex157
- mutex attributes
 - extended160
- mutex initialization171
- mv228, 275
- name space88
- name space pollution87-88
- NAME_MAX223, 286
- NaN arguments41
- nanosleep()136, 138, 140, 273
- native implementation23
- network byte order36
- newgrp276
- news, rationale for omission262
- NGROUPS_MAX4, 25, 223, 280, 286

Index

nice.....	275	pathname resolution.....	38
nice value.....	23	pathname, incomplete.....	18
nice().....	130	PATH_MAX.....	93, 223, 286
nl_langinfo().....	273	pattern matching	
nm.....	277	multiple character.....	253
noclobber option.....	242	notation.....	252
nohup.....	275	single character.....	252
non-canonical mode input processing.....	69	patterns	
non-printable.....	31	filename expansion.....	253
normative references.....	5, 79, 219	pax.....	275-277
NQS.....	255	pcat, rationale for omission.....	262
obsolescent.....	6	pclose().....	277
obsolescent, rationale.....	6	Pending error.....	179
od.....	275, 277	per-thread errno.....	94
open.....	272	performance enhancements.....	268
open file description.....	23	pg, rationale for omission.....	262
open file descriptors.....	243	PID_MAX.....	205
open().....	18, 68, 120-122, 228	pipe.....	19, 24, 272
opendir.....	272	pipe().....	101
openlog().....	276	pipelines.....	246
OPEN_MAX.....	5, 223, 286-287	PIPE_BUF.....	223, 286
option definitions.....	255	popen().....	274, 277-278
optional behavior.....	288	portability.....	8, 85, 221
options.....	179	portability codes.....	8, 85, 221
OR lists.....	248	portable character set.....	42
orphaned process group.....	23, 102	portable filename character set.....	24
output device.....	65	positional parameters.....	234
output file descriptor		POSIX locale.....	45
duplication.....	243	POSIX.1 symbols.....	87
output mode.....	71	POSIX.13.....	123
output processing.....	70	POSIX2_BC_BASE_MAX.....	279
overrun conditions.....	204	POSIX2_BC_DIM_MAX.....	279
overrun in dumping trace streams.....	204	POSIX2_BC_SCALE_MAX.....	279
overrun in trace streams.....	204	POSIX2_BC_STRING_MAX.....	279
pack, rationale for omission.....	262	POSIX2_CHAR_TERM.....	279
page.....	24, 119, 123	POSIX2_COLL_WEIGHTS_MAX.....	279
PAGESIZE.....	119, 163, 286	POSIX2_C_BIND.....	224, 278
parallel I/O.....	164	POSIX2_C_DEV.....	224, 278
parameter expansion.....	238	POSIX2_EXPR_NEST_MAX.....	279
parameters.....	70, 234	POSIX2_FORT_DEV.....	224, 278
parameters, positional.....	234	POSIX2_FORT_RUN.....	224, 278
parameters, special.....	234	POSIX2_LINE_MAX.....	279
parent directory.....	24	POSIX2_LOCALEDEF.....	224, 276, 278
passwd file.....	24	POSIX2_PBS.....	279
passwd, rationale for omission.....	262	POSIX2_PBS_ACCOUNTING.....	279
paste.....	275	POSIX2_PBS_CHECKPOINT.....	279
patch.....	275-277	POSIX2_PBS_LOCATE.....	279
PATH.....	57	POSIX2_PBS_MESSAGE.....	279
pathchk.....	275	POSIX2_PBS_TRACK.....	279
pathconf().....	18, 222, 271-2	POSIX2_RE_DUP_MAX.....	280
pathname expansion.....	241	POSIX2_SW_DEV.....	224, 278

POSIX2_SYMLINKS	223	POSIX_TRACE_LOOP	204
POSIX2_UPE	224, 278-279	posix_typed_mem_get_info()	123
POSIX2_VERSION	280	posix_typed_mem_open()	123
POSIX_ALLOC_SIZE_MIN	106	POSIX_USER_GROUPS	294
POSIX_C_LANG_JUMP	291	POSIX_USER_GROUPS_R	294
POSIX_C_LANG_MATH	291	POSIX_VERSION	286
POSIX_C_LANG_SUPPORT	292	POSIX_WIDE_CHAR_DEVICE_IO	294
POSIX_C_LANG_SUPPORT_R	292	post-mortem filtering of trace event types	201
POSIX_C_LANG_WIDE_CHAR	292	pr	275-276
POSIX_C_LIB_EXT	292	pread()	164
POSIX_DEVICE_IO	292	printf	274-275
POSIX_DEVICE_SPECIFIC	292	printing	270
POSIX_DEVICE_SPECIFIC_R	292	privilege	33
posix_fadvise()	106	process group	67
POSIX_FADV_DONTNEED	106	process group ID	19, 67
POSIX_FADV_NOREUSE	106	process group lifetime	67
POSIX_FADV_RANDOM	106	process group, orphaned	23, 102
POSIX_FADV_SEQUENTIAL	106	process groups, concepts in job control	19
POSIX_FADV_WILLNEED	106	process ID reuse	39
POSIX_FD_MGMT	292	process ID, rationale	205
POSIX_FIFO	292	process management	268, 271
POSIX_FILE_ATTRIBUTES	292	process scheduling	129, 271
POSIX_FILE_LOCKING	292	prof, rationale for omission	262
POSIX_FILE_SYSTEM	293	profiling	277
POSIX_FILE_SYSTEM_EXT	293	programming manipulation	194
POSIX_FILE_SYSTEM_R	293	prompting	236
POSIX_JOB_CONTROL	293	protocols	178
posix_madvise()	106	ps	275-276
POSIX_MADV_DONTNEED	106	pthread	202
POSIX_MADV_RANDOM	106	pthread_attr_getguardsize()	164
POSIX_MADV_SEQUENTIAL	106	pthread_attr_setguardsize()	164
POSIX_MADV_WILLNEED	106	PTHREAD_BARRIER_SERIAL_THREAD	155
posix_mem_offset()	123-124	pthread_barrier_wait()	155, 175
POSIX_MULTI_PROCESS	293	pthread_cond_init()	152
POSIX_NETWORKING	293	pthread_cond_timedwait()	94, 139, 161, 280
POSIX_PIPE	293	pthread_cond_wait()	94, 111, 161
POSIX_REC_INCR_XFER_SIZE	107	pthread_create()	152-153
POSIX_REC_MAX_XFER_SIZE	107	PTHREAD_CREATE_DETACHED	173
POSIX_REC_MIN_XFER_SIZE	107	PTHREAD_DESTRUCTOR_ITERATIONS	287
POSIX_REC_XFER_ALIGN	106	pthread_detach()	173
POSIX_REGEXP	293	pthread_getconcurrency()	163
POSIX_SHELL_FUNC	293	pthread_getcpuclockid()	142-143
POSIX_SIGNALS	293	pthread_join()	94, 173
POSIX_SIGNAL_JUMP	293	PTHREAD_KEYS_MAX	287
POSIX_SINGLE_PROCESS	293	pthread_key_create()	154
posix_spawn()	207, 271	pthread_key_gettype()	161
posix_spawnnp()	207, 271	pthread_mutexattr_settype()	161
POSIX_SYMBOLIC_LINKS	293	PTHREAD_MUTEX_DEFAULT	160
POSIX_SYSTEM_DATABASE	293	PTHREAD_MUTEX_ERRORCHECK	160
POSIX_SYSTEM_DATABASE_R	293	pthread_mutex_init()	152
posix_trace_eventid_open()	199	pthread_mutex_lock()	94, 160, 174

Index

- PTHREAD_MUTEX_NORMAL160
- PTHREAD_MUTEX_RECURSIVE160
- pthread_mutex_timedlock()140
- pthread_mutex_trylock()160
- pthread_mutex_unlock()160
- PTHREAD_PROCESS_PRIVATE162
- PTHREAD_PROCESS_SHARED162
- pthread_rwlockattr_destroy()162
- pthread_rwlockattr_getpshared()162
- pthread_rwlockattr_init()162
- pthread_rwlockattr_setpshared()162
- pthread_rwlock_init()162
- PTHREAD_RWLOCK_INITIALIZER162
- pthread_rwlock_rdlock()162
- pthread_rwlock_t162
- pthread_rwlock_tryrdlock()162
- pthread_rwlock_trywrlock()162
- pthread_rwlock_unlock()163, 176
- pthread_rwlock_wrlock()162
- pthread_self()154
- pthread_setconcurrency()163
- pthread_setprio()171
- pthread_setschedparam()171
- pthread_setspecific()154
- pthread_spin_lock()156, 175
- pthread_spin_trylock()156
- PTHREAD_STACK_MIN287
- PTHREAD_THREADS_MAX287
- putc()165, 272
- putchar272
- pwd275
- pwrite()164
- queuing of waiting threads176
- quote removal241
- quoting231
- rand()174
- RCS, rationale for omission262
- RE
 - grammar65
- RE bracket expression60
- read230, 272, 274
- read lock162
- read()19, 68-69, 92, 101-102, 104
 -113-115, 119-120, 164, 174, 206
- read-write attributes161
- read-write locks161
- readdir272
- reading an active trace stream204
- reading data69
- readlink()26, 272
- realpath272
- realtime106
- realtime signal delivery98
- realtime signal generation98
- realtime signals111
- red, rationale for omission262
- redirect input243
- redirect output243
- redirection241
- references5, 79, 219
- regcomp()277
- regerror()277
- regexec()277
- regfree()277
- regular expression58
 - definitions58
 - general requirements59
 - grammar64
- regular file25
- rejected utilities260
- remove()26
- rename()26, 272
- renice275
- replenishment period132
- reserved words233
- rewinddir272
- RE_DUP_MAX223
- rm228, 275
- rmdir272, 275
- rmdir()26, 93
- root directory25, 38-39
- root file system25
- root of a file system25
- routing178
- rsh, rationale for omission262
- RTSIG_MAX287
- samefile()205
- SA_NOCLDSTOP20
- SA_SIGINFO99-100
- scheduling allocation domain169
- scheduling contention scope169-170
- scheduling documentation170
- scheduling policy39
- SCHED_FIFO130-131, 169, 176, 273
- SCHED_OTHER130
- SCHED_RR130, 169, 176, 273
- SCHED_SPORADIC273
- scope3, 79, 219
- sdb, rationale for omission262
- sdiff, rationale for omission262
- seconds since the Epoch39
- security considerations13, 17, 21, 32-33, 68

- security, monolithic privileges 13
- sed 275-276
- sem*() 105
- semaphore 40, 109, 273
- semctl() 105
- semget() 105
- semop() 105
- sem_init() 109
- SEM_NSEMS_MAX 287
- sem_open() 109
- sem_timedwait() 140
- sem_trywait() 94, 111
- SEM_VALUE_MAX 287
- sem_wait() 94, 111
- sequential lists 248
- session 20, 24, 67
- set 236
- setgid() 25
- setgrent() 32
- setjmp() 274
- setlocale() 273
- setlogmask() 276
- setpgid() 19-20, 67
- setpriority() 130
- setpwent() 32
- setrlimit() 229
- setsid() 67
- setuid() 25
- sh 276, 283
- shall 6
- shall, rationale 6
- shar, rationale for omission 262
- shared memory 121
- shell 18-21, 67-68, 95, 101-102
- SHELL 57
- shell commands 244
- shell errors 243
- shell execution environment 234, 252
- shell grammar 251
 - lexical conventions 252
 - rules 252
- shell variables 235
- shell, job control 19, 95, 102
- shl, rationale for omission 262
- shm*() 105
- shmctl() 105
- shmdt() 105
- shm_open() 120-123
- shm_unlink() 121-123
- should 6
- should, rationale 6
- SIGABRT 30, 95
- sigaction() 97, 99
- SIGBUS 30, 95
- SIGCHLD 20, 97, 101
- SIGCLD 101
- SIGCONT 20, 98, 101-102
- SIGEMT 95
- SIGEV_NONE 98
- SIGEV_NOTIFY 98
- SIGEV_SIGNAL 98
- SIGFPE 30, 95, 97, 100
- SIGHUP 102
- SIGILL 30, 95
- SIGINT 21, 167
- SIGIOT 95
- SIGKILL 95, 98, 102
- siglongjmp() 92, 103, 274
- signal 25
- signal acceptance 96
- signal actions 101
- signal concepts 95
- signal delivery 96
- signal generation 96
- signal names 95
- signal() 95, 97
- signals 179, 252
- SIGPIPE 30, 93
- sigprocmask() 97
- SIGRTMAX 99-100
- SIGRTMIN 99-100
- SIGSEGV 30, 95, 100
- sigsetjmp() 274
- sigset_t 95
- SIGSTOP 102
- sigsuspend() 101, 104
- SIGSYS 95
- SIGTERM 95
- sigtimedwait() 94, 114, 139
- SIGTRAP 95
- SIGTSTP 21, 102
- SIGTTIN 20, 68, 102
- SIGTTOU 20, 68, 102
- SIGUSR1 95
- SIGUSR2 95
- sigwait() 94, 174
- sigwaitinfo() 94, 114
- sigwait_multiple() 97
- SIG_DFL 97-98, 101
- SIG_IGN 20, 97-98, 101, 103
- simple commands 244
- single-quotes 231

Index

size, rationale for omission.....	262	symbolic link	26
SI_USER	99-100	symbolic name	7
sleep	271, 274	symbols.....	87
sleep()	103-104, 273	SYMLOOP_MAX.....	92
socket I/O mode.....	178	synchronized I/O.....	114, 272
socket out-of-band data state.....	179	data integrity completion.....	31, 114
socket owner.....	179	file integrity completion	31, 114
socket queue limit	179	synchronously-generated signal	30
socket receive queue	179	sysconf()	18, 117, 119, 121, 167, 222, 271-272
socket types	178	syslog()	276
sockets.....	178	system call.....	31
Internet Protocols	179	system documentation	6
IPv4.....	179	system environment	270, 276
IPv6.....	179	System III	24, 205
local UNIX connection	179	system interfaces	207
software development	270, 277	system reboot	31
sort	275-276	System V.....	16, 21, 96-97, 101
spawn example	207	system().....	274, 277-278
special built-in.....	249	tabs.....	275
special built-in utilities	254	tail	275
special characters.....	70	talk.....	274, 276
special control character	71	tar, rationale for omission.....	263
special parameters.....	234	tcgetattr()	20
specific implementation.....	18	tcgetpgrp().....	20, 67
spell, rationale for omission	263	tcsetattr().....	20, 66
spin locks.....	156-157	tcsetpgrp().....	19-20
split.....	275	tee.....	274
sporadic server scheduling policy	132	TERM	57
SSIZE_MAX.....	206, 287	terminal access control.....	68
SS_REPL_MAX	135	terminal device file.....	67
standard I/O streams	104	closing.....	70
stat().....	15, 120, 228, 271-2	terminal type	65
state-dependent character encoding	42	terminology	5, 85, 220
statvfs()	229	termios structure.....	70
STREAMS.....	104	test.....	274-275
STREAM_MAX.....	287	TeX.....	275
strings.....	277	text file.....	31
strip.....	277	thread	31
structures, additions to.....	88	thread cancelability states	174
stty	57, 276	thread cancelability type.....	174
su, rationale for omission.....	263	thread cancelation.....	172, 174
subprofiling.....	11	thread cancelation points.....	174
subprofiling option groups	291	thread concurrency level.....	163
subshells	20	thread creation attributes.....	151
successfully completed	31	thread ID	32, 167
sum	228	thread interactions	178
sum, rationale for omission.....	263	thread mutex	168
superuser.....	13, 25, 33, 236, 260	thread read-write lock	176
supplementary group ID	25	thread scheduling.....	168
supplementary groups	33	thread stack guard size.....	163
symbolic constant.....	7	thread-safe.....	32

- thread-safe function32
- thread-safety.....40, 164
- thread-safety, rationale.....40
- thread-specific data.....153
- threads151
 - implementation models.....153
- tilde expansion.....237
- time.....274, 277
- time()92
- timeouts.....144
- timers135
- TIMER_ABSTIME.....136-137
- TIMER_MAX.....287
- timer_settime()136-137
- times()92, 143, 271
- timestamp clock.....203
- time_t40
- token recognition.....232
- TOSTOP.....20
- touch228, 275
- tput276
- tr.275-276
- trace analyzer192
- trace event type-filtering.....201
- trace event types.....201
- trace examples.....190
- trace model185
- trace operation control190
- trace storage.....189
- trace stream attribute.....195
- trace stream states.....188
- tracing.....40, 180
- tracing all processes.....188
- tracing, detailed objectives.....181
- triggering.....203
- troff275
- trojan horse13
- true230, 274
- tsort, rationale for omission263
- tty.....276
- ttyname()165
- TTY_NAME_MAX.....287
- typed memory.....123
- TZ.....57
- TZNAME_MAX.....287
- tzset()273
- ualarm.....271
- UID_MAX206
- uid_t.....32
- ULONG_MAX222
- umask230, 276
- umask()271
- umount()22
- unalias.....229, 274
- uname276
- uname()271
- unbounded priority inversion171
- undefined.....6
- undefined, rationale6
- unexpand.....275
- uniq.....275-276
- unlink272
- unlink()26, 93, 120-121, 123
- unpack, rationale for omission263
- unsafe functions.....102
- unspecified6
- unspecified, rationale.....6
- until loop249
- user database32
- user database access.....32
- user requirements.....267
- usleep271
- utility41
- utility argument syntax.....72
- utility conventions.....72
- utility description defaults225
- utility limits222
- utility syntax guidelines.....73
- utime()26, 272
- uucp.....276
- uudecode275-276
- uencode275-276
- variable assignment41
- variables234
- VEOF71
- VEOL71
- Version 7.....38, 230
- vhangup()21
- vi274, 276
- virtual processor33
- VMIN71
- VTIME.....71
- wait.....230, 274
- wait()92, 95, 101-102, 104, 271
- waitpid()20, 24, 102, 205, 271
- wall, rationale for omission.....263
- wc.....275
- WERASE.....69
- while loop.....249
- who.....276
- wide-character codes42
- word expansions.....236

Index

wordexp()	277
wordfree()	277
write	272, 274, 276
write lock	162
write()	19-20, 68, 92, 101-102 104, 113-115, 119-120, 206
writing data	70
WUNTRACED	20
xargs	274
XSI	33
XSI IPC	105
XSI supported functions	158
XSI threads extensions	159
XSI_C_LANG_SUPPORT	294
XSI_DBM	294
XSI_DEVICE_IO	294
XSI_DEVICE_SPECIFIC	294
XSI_DYNAMIC_LINKING	294
XSI_FD_MGMT	294
XSI_FILE_SYSTEM	294
XSI_I18N	294
XSI_IPC	294
XSI_JOB_CONTROL	294
XSI_JUMP	294
XSI_MATH	294
XSI_MULTI_PROCESS	294
XSI_SIGNALS	295
XSI_SINGLE_PROCESS	295
XSI_SYSTEM_DATABASE	295
XSI_SYSTEM_LOGGING	295
XSI_THREADS_EXT	295
XSI_THREAD_MUTEX_EXT	295
XSI_TIMERS	295
XSI_USER_GROUPS	295
XSI_WIDE_CHAR	295
yacc	277-278

