

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

Draft Technical Standard: Rationale, Issue 7

Prepared by the Austin Group (www.opengroup.org/austin)

Copyright © 2006 The Institute of Electrical & Electronics Engineers, Inc.
3 Park Avenue, New York, NY 10016-5997, USA

Copyright © 2006 The Open Group
Thames Tower, Station Road, Reading, Berkshire RG1 1LX, UK

All rights reserved.

Except as permitted below, no part of this publication may be reproduced, stored in a retrieval system, or transmitted, in any form or by any means, electronic, mechanical, photocopying, recording, or otherwise, without the prior permission of the copyright owners. This is an unapproved draft, subject to change. Permission is hereby granted for Austin Group participants to reproduce this document for purposes of IEEE, The Open Group, and JTC1 standardization activities. Other entities seeking permission to reproduce this document for standardization purposes or other activities must contact the copyright owners for an appropriate license. Use of information contained within this unapproved draft is at your own risk.

Portions of this document 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.

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. 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, which contains 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)

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 www.opengroup.org/austin/bugreport.html.

Contents

1	Part	A	Base Definitions	i
2	Appendix	A	Rationale for Base Definitions.....	3
3		A.1	Introduction	3
4		A.1.1	Scope.....	3
5		A.1.2	Conformance	6
6		A.1.3	Normative References.....	6
7		A.1.4	Terminology.....	6
8		A.1.5	Portability	8
9		A.1.5.1	Codes	8
10		A.1.5.2	Margin Code Notation.....	9
11		A.2	Conformance.....	9
12		A.2.1	Implementation Conformance.....	9
13		A.2.1.1	Requirements	9
14		A.2.1.2	Documentation.....	9
15		A.2.1.3	POSIX Conformance	10
16		A.2.1.4	XSI Conformance.....	11
17		A.2.1.5	Option Groups	11
18		A.2.1.6	Options.....	12
19		A.2.2	Application Conformance.....	12
20		A.2.2.1	Strictly Conforming POSIX Application	12
21		A.2.2.2	Conforming POSIX Application.....	12
22		A.2.2.3	Conforming POSIX Application Using Extensions	12
23		A.2.2.4	Strictly Conforming XSI Application.....	12
24		A.2.2.5	Conforming XSI Application Using Extensions.....	12
25		A.2.3	Language-Dependent Services for the C Programming Language	13
26		A.2.4	Other Language-Related Specifications	13
27		A.3	Definitions	13
28		A.4	General Concepts.....	33
29		A.4.1	Concurrent Execution	33
30		A.4.2	Directory Protection	33
31		A.4.3	Extended Security Controls.....	34
32		A.4.4	File Access Permissions	34
33		A.4.5	File Hierarchy	34
34		A.4.6	Filenames	34
35		A.4.7	Filename Portability	36
36		A.4.8	File Times Update	36
37		A.4.9	Host and Network Byte Order	36
38		A.4.10	Measurement of Execution Time.....	36
39		A.4.11	Memory Synchronization.....	37
40		A.4.12	Pathname Resolution	38
41		A.4.13	Process ID Reuse	39
42		A.4.14	Scheduling Policy	40

43	A.4.15	Seconds Since the Epoch.....	40
44	A.4.16	Semaphore	41
45	A.4.17	Thread-Safety	41
46	A.4.18	Tracing	41
47	A.4.19	Treatment of Error Conditions for Mathematical Functions	41
48	A.4.20	Treatment of NaN Arguments for Mathematical Functions	41
49	A.4.21	Utility	41
50	A.4.22	Variable Assignment	41
51	A.5	File Format Notation	42
52	A.6	Character Set	42
53	A.6.1	Portable Character Set.....	42
54	A.6.2	Character Encoding.....	43
55	A.6.3	C Language Wide-Character Codes.....	43
56	A.6.4	Character Set Description File	43
57	A.6.4.1	State-Dependent Character Encodings.....	43
58	A.7	Locale.....	45
59	A.7.1	General	45
60	A.7.2	POSIX Locale	46
61	A.7.3	Locale Definition.....	46
62	A.7.3.1	LC_CTYPE	47
63	A.7.3.2	LC_COLLATE	48
64	A.7.3.3	LC_MONETARY.....	50
65	A.7.3.4	LC_NUMERIC.....	51
66	A.7.3.5	LC_TIME.....	51
67	A.7.3.6	LC_MESSAGES.....	52
68	A.7.4	Locale Definition Grammar	52
69	A.7.4.1	Locale Lexical Conventions.....	52
70	A.7.4.2	Locale Grammar	52
71	A.7.5	Locale Definition Example	52
72	A.8	Environment Variables.....	56
73	A.8.1	Environment Variable Definition	56
74	A.8.2	Internationalization Variables.....	56
75	A.8.3	Other Environment Variables	57
76	A.9	Regular Expressions	58
77	A.9.1	Regular Expression Definitions	59
78	A.9.2	Regular Expression General Requirements	60
79	A.9.3	Basic Regular Expressions	61
80	A.9.3.1	BREs Matching a Single Character or Collating Element	61
81	A.9.3.2	BRE Ordinary Characters	61
82	A.9.3.3	BRE Special Characters	61
83	A.9.3.4	Periods in BREs	61
84	A.9.3.5	RE Bracket Expression	61
85	A.9.3.6	BREs Matching Multiple Characters.....	62
86	A.9.3.7	BRE Precedence.....	63
87	A.9.3.8	BRE Expression Anchoring	63
88	A.9.4	Extended Regular Expressions	63
89	A.9.4.1	EREs Matching a Single Character or Collating Element	64
90	A.9.4.2	ERE Ordinary Characters	64
91	A.9.4.3	ERE Special Characters	64
92	A.9.4.4	Periods in EREs	64
93	A.9.4.5	ERE Bracket Expression.....	64
94	A.9.4.6	EREs Matching Multiple Characters.....	64

Contents

95	A.9.4.7	ERE Alternation	64
96	A.9.4.8	ERE Precedence.....	64
97	A.9.4.9	ERE Expression Anchoring	64
98	A.9.5	Regular Expression Grammar	64
99	A.9.5.1	BRE/ERE Grammar Lexical Conventions	65
100	A.9.5.2	RE and Bracket Expression Grammar	65
101	A.9.5.3	ERE Grammar	65
102	A.10	Directory Structure and Devices	65
103	A.10.1	Directory Structure and Files	65
104	A.10.2	Output Devices and Terminal Types	66
105	A.11	General Terminal Interface.....	66
106	A.11.1	Interface Characteristics	67
107	A.11.1.1	Opening a Terminal Device File	67
108	A.11.1.2	Process Groups.....	67
109	A.11.1.3	The Controlling Terminal	68
110	A.11.1.4	Terminal Access Control.....	68
111	A.11.1.5	Input Processing and Reading Data	69
112	A.11.1.6	Canonical Mode Input Processing	69
113	A.11.1.7	Non-Canonical Mode Input Processing	69
114	A.11.1.8	Writing Data and Output Processing.....	70
115	A.11.1.9	Special Characters.....	70
116	A.11.1.10	Modem Disconnect.....	70
117	A.11.1.11	Closing a Terminal Device File	70
118	A.11.2	Parameters that Can be Set.....	70
119	A.11.2.1	The termios Structure.....	70
120	A.11.2.2	Input Modes	70
121	A.11.2.3	Output Modes	70
122	A.11.2.4	Control Modes.....	71
123	A.11.2.5	Local Modes.....	71
124	A.11.2.6	Special Control Characters	71
125	A.12	Utility Conventions	71
126	A.12.1	Utility Argument Syntax	71
127	A.12.2	Utility Syntax Guidelines	72
128	A.13	Headers	75
129	A.13.1	Format of Entries	75
130	A.13.2	Removed Headers in Issue 7	75
131	Part B	System Interfaces	lxxvii
132	Appendix B	Rationale for System Interfaces	79
133	B.1	Introduction.....	79
134	B.1.1	Scope.....	79
135	B.1.2	Conformance	79
136	B.1.3	Normative References.....	79
137	B.1.4	Change History	79
138	B.1.5	Terminology.....	82
139	B.1.6	Definitions	83
140	B.1.7	Relationship to Other Formal Standards.....	83
141	B.1.8	Portability	83
142	B.1.8.1	Codes	83
143	B.1.9	Format of Entries	83

144	B.2	General Information.....	83
145	B.2.1	Use and Implementation of Functions	83
146	B.2.2	The Compilation Environment.....	84
147	B.2.2.1	POSIX.1 Symbols	84
148	B.2.2.2	The Name Space	85
149	B.2.3	Error Numbers	89
150	B.2.3.1	Additional Error Numbers.....	93
151	B.2.4	Signal Concepts.....	93
152	B.2.4.1	Signal Generation and Delivery	94
153	B.2.4.2	Realtime Signal Generation and Delivery.....	96
154	B.2.4.3	Signal Actions.....	99
155	B.2.4.4	Signal Effects on Other Functions	102
156	B.2.5	Standard I/O Streams	102
157	B.2.5.1	Interaction of File Descriptors and Standard I/O Streams.....	102
158	B.2.5.2	Stream Orientation and Encoding Rules.....	102
159	B.2.6	STREAMS	102
160	B.2.6.1	Accessing STREAMS.....	103
161	B.2.7	XSI Interprocess Communication.....	103
162	B.2.7.1	IPC General Information	103
163	B.2.8	Realtime	104
164	B.2.8.1	Realtime Signals.....	109
165	B.2.8.2	Asynchronous I/O	111
166	B.2.8.3	Memory Management.....	113
167	B.2.8.4	Process Scheduling	127
168	B.2.8.5	Clocks and Timers	133
169	B.2.9	Threads.....	149
170	B.2.9.1	Thread-Safety	163
171	B.2.9.2	Thread IDs	166
172	B.2.9.3	Thread Mutexes	166
173	B.2.9.4	Thread Scheduling.....	166
174	B.2.9.5	Thread Cancellation	171
175	B.2.9.6	Thread Read-Write Locks	175
176	B.2.9.7	Thread Interactions with Regular File Operations	176
177	B.2.9.8	Use of Application-Managed Thread Stacks	176
178	B.2.10	Sockets.....	177
179	B.2.10.1	Address Families	177
180	B.2.10.2	Addressing.....	177
181	B.2.10.3	Protocols.....	177
182	B.2.10.4	Routing.....	177
183	B.2.10.5	Interfaces.....	177
184	B.2.10.6	Socket Types	177
185	B.2.10.7	Socket I/O Mode	177
186	B.2.10.8	Socket Owner	177
187	B.2.10.9	Socket Queue Limits	177
188	B.2.10.10	Pending Error	177
189	B.2.10.11	Socket Receive Queue	178
190	B.2.10.12	Socket Out-of-Band Data State	178
191	B.2.10.13	Connection Indication Queue	178
192	B.2.10.14	Signals	178
193	B.2.10.15	Asynchronous Errors	178
194	B.2.10.16	Use of Options.....	178
195	B.2.10.17	Use of Sockets for Local UNIX Connections.....	178

Contents

196	B.2.10.18	Use of Sockets over Internet Protocols	178
197	B.2.10.19	Use of Sockets over Internet Protocols Based on IPv4	178
198	B.2.10.20	Use of Sockets over Internet Protocols Based on IPv6	178
199	B.2.11	Tracing	178
200	B.2.11.1	Objectives	178
201	B.2.11.2	Trace Model	183
202	B.2.11.3	Trace Programming Examples	188
203	B.2.11.4	Rationale on Trace for Debugging	196
204	B.2.11.5	Rationale on Trace Event Type Name Space	196
205	B.2.11.6	Rationale on Trace Events Type Filtering	198
206	B.2.11.7	Tracing, pthread API	200
207	B.2.11.8	Rationale on Triggering	201
208	B.2.11.9	Rationale on Timestamp Clock	201
209	B.2.11.10	Rationale on Different Overrun Conditions	202
210	B.2.12	Data Types	202
211	B.2.12.1	Defined Types	202
212	B.2.12.2	The char Type	204
213	B.2.12.3	Pointer Types	204
214	B.3	System Interfaces	205
215	B.3.1	System Interfaces Removed in this Revision	205
216	B.3.1.1	bcmp	205
217	B.3.1.2	bcopy	205
218	B.3.1.3	bsd_signal	205
219	B.3.1.4	bzero	205
220	B.3.1.5	ecvt, fcvt, gcvt	205
221	B.3.1.6	ftime	205
222	B.3.1.7	getcontext, makecontext, swapcontext	206
223	B.3.1.8	gethostbyaddr, gethostbyname	206
224	B.3.1.9	getwd	206
225	B.3.1.10	h_errno	206
226	B.3.1.11	index	206
227	B.3.1.12	makecontext	206
228	B.3.1.13	mktemp	206
229	B.3.1.14	pthread_attr_getstackaddr, pthread_attr_setstackaddr	206
230	B.3.1.15	rindex	206
231	B.3.1.16	scalb	207
232	B.3.1.17	ualarm	207
233	B.3.1.18	usleep	207
234	B.3.1.19	vfork	207
235	B.3.1.20	wcs wcs	207
236	B.3.2	Examples for Spawn	207
237	Part C	Shell and Utilities	ccxvii
238	Appendix C	Rationale for Shell and Utilities	219
239	C.1	Introduction	219
240	C.1.1	Scope	219
241	C.1.2	Conformance	219
242	C.1.3	Normative References	219
243	C.1.4	Change History	219
244	C.1.5	Terminology	220

245	C.1.6	Definitions	220
246	C.1.7	Relationship to Other Documents.....	220
247	C.1.7.1	System Interfaces	220
248	C.1.7.2	Concepts Derived from the ISO C Standard.....	220
249	C.1.8	Portability	221
250	C.1.8.1	Codes	221
251	C.1.9	Utility Limits	221
252	C.1.10	Grammar Conventions	224
253	C.1.11	Utility Description Defaults	224
254	C.1.12	Considerations for Utilities in Support of Files of Arbitrary Size	228
255	C.1.13	Built-In Utilities.....	228
256	C.2	Shell Command Language.....	230
257	C.2.1	Shell Introduction	230
258	C.2.2	Quoting	230
259	C.2.2.1	Escape Character (Backslash).....	230
260	C.2.2.2	Single-Quotes	230
261	C.2.2.3	Double-Quotes	230
262	C.2.3	Token Recognition	232
263	C.2.3.1	Alias Substitution	232
264	C.2.4	Reserved Words	233
265	C.2.5	Parameters and Variables	233
266	C.2.5.1	Positional Parameters.....	233
267	C.2.5.2	Special Parameters.....	233
268	C.2.5.3	Shell Variables	234
269	C.2.6	Word Expansions	236
270	C.2.6.1	Tilde Expansion.....	236
271	C.2.6.2	Parameter Expansion	237
272	C.2.6.3	Command Substitution.....	237
273	C.2.6.4	Arithmetic Expansion	238
274	C.2.6.5	Field Splitting.....	241
275	C.2.6.6	Pathname Expansion.....	241
276	C.2.6.7	Quote Removal	241
277	C.2.7	Redirection.....	241
278	C.2.7.1	Redirecting Input.....	243
279	C.2.7.2	Redirecting Output.....	243
280	C.2.7.3	Appending Redirected Output.....	243
281	C.2.7.4	Here-Document.....	243
282	C.2.7.5	Duplicating an Input File Descriptor.....	243
283	C.2.7.6	Duplicating an Output File Descriptor.....	243
284	C.2.7.7	Open File Descriptors for Reading and Writing	243
285	C.2.8	Exit Status and Errors.....	243
286	C.2.8.1	Consequences of Shell Errors.....	243
287	C.2.8.2	Exit Status for Commands.....	243
288	C.2.9	Shell Commands.....	244
289	C.2.9.1	Simple Commands	244
290	C.2.9.2	Pipelines	246
291	C.2.9.3	Lists.....	247
292	C.2.9.4	Compound Commands	248
293	C.2.9.5	Function Definition Command.....	249
294	C.2.10	Shell Grammar	251
295	C.2.10.1	Shell Grammar Lexical Conventions	252
296	C.2.10.2	Shell Grammar Rules	252

Contents

297	C.2.11	Signals and Error Handling	252
298	C.2.12	Shell Execution Environment.....	252
299	C.2.13	Pattern Matching Notation.....	252
300	C.2.13.1	Patterns Matching a Single Character	252
301	C.2.13.2	Patterns Matching Multiple Characters	253
302	C.2.13.3	Patterns Used for Filename Expansion	253
303	C.2.14	Special Built-In Utilities	254
304	C.3	Batch Environment Services and Utilities.....	254
305	C.3.1	Batch General Concepts.....	257
306	C.3.2	Batch Services.....	259
307	C.3.3	Common Behavior for Batch Environment Utilities	260
308	C.4	Utilities	260
309	Part D	Portability Considerations	cclxv
310	Appendix D	Portability Considerations (Informative).....	267
311	D.1	User Requirements	267
312	D.1.1	Configuration Interrogation.....	268
313	D.1.2	Process Management.....	268
314	D.1.3	Access to Data	268
315	D.1.4	Access to the Environment.....	268
316	D.1.5	Access to Determinism and Performance Enhancements	268
317	D.1.6	Operating System-Dependent Profile.....	268
318	D.1.7	I/O Interaction	268
319	D.1.8	Internationalization Interaction.....	269
320	D.1.9	C-Language Extensions	269
321	D.1.10	Command Language.....	269
322	D.1.11	Interactive Facilities.....	269
323	D.1.12	Accomplish Multiple Tasks Simultaneously	269
324	D.1.13	Complex Data Manipulation	269
325	D.1.14	File Hierarchy Manipulation.....	269
326	D.1.15	Locale Configuration.....	269
327	D.1.16	Inter-User Communication	270
328	D.1.17	System Environment.....	270
329	D.1.18	Printing.....	270
330	D.1.19	Software Development	270
331	D.2	Portability Capabilities	270
332	D.2.1	Configuration Interrogation.....	271
333	D.2.2	Process Management.....	271
334	D.2.3	Access to Data	272
335	D.2.4	Access to the Environment.....	273
336	D.2.5	Bounded (Realtime) Response.....	273
337	D.2.6	Operating System-Dependent Profile.....	273
338	D.2.7	I/O Interaction	274
339	D.2.8	Internationalization Interaction.....	274
340	D.2.9	C-Language Extensions	274
341	D.2.10	Command Language.....	274
342	D.2.11	Interactive Facilities.....	275
343	D.2.12	Accomplish Multiple Tasks Simultaneously	275
344	D.2.13	Complex Data Manipulation	275
345	D.2.14	File Hierarchy Manipulation.....	276

346	D.2.15	Locale Configuration.....	276
347	D.2.16	Inter-User Communication	276
348	D.2.17	System Environment	277
349	D.2.18	Printing.....	277
350	D.2.19	Software Development	277
351	D.2.20	Future Growth.....	277
352	D.3	Profiling Considerations.....	278
353	D.3.1	Configuration Options.....	278
354	D.3.2	Configuration Options (Shell and Utilities).....	278
355	D.3.3	Configurable Limits	280
356	D.3.4	Configuration Options (System Interfaces)	280
357	D.3.5	Configurable Limits	285
358	D.3.6	Optional Behavior.....	287
359	Part E	Subprofiling Considerations.....	cclxxxix
360	Appendix E	Subprofiling Considerations (Informative).....	291
361	E.1	Subprofiling Option Groups	291
362		Index.....	297
363	List of Figures		
364	B-1	Example of a System with Typed Memory	122
365	B-2	Trace System Overview: for Offline Analysis	184
366	B-3	Trace System Overview: for Online Analysis	185
367	B-4	Trace System Overview: States of a Trace Stream	186
368	B-5	Trace Another Process	196
369	B-6	Trace Name Space Overview: With Third-Party Library	197
370	List of Tables		
371	A-1	Historical Practice for Symbolic Links.....	30

372

Foreword

373

Structure of the Standard

374

Notes to Reviewers

375

This section with side shading will not appear in the final copy. - Ed.

376

This section will be completed in a later draft.

DRAFT

Introduction

378 **Note:** This introduction is not part of IEEE Std 1003.1-200x, Standard for Information Technology —
379 Portable Operating System Interface (POSIX).

380 This draft standard was developed, and is maintained, by a joint working group of members of
381 the IEEE Portable Applications Standards Committee, members of The Open Group, and
382 members of ISO/IEC Joint Technical Committee 1. This joint working group is known as the
383 Austin Group.¹

384 The Austin Group arose out of discussions amongst the parties which started in early 1998,
385 leading to an initial meeting and formation of the group in September 1998. The purpose of the
386 Austin Group is to develop and maintain the core open systems interfaces that are the POSIX[®]
387 1003.1 (and former 1003.2) standards, ISO/IEC 9945 Parts 1 to 4, and the core of the Single UNIX
388 Specification.

389 The approach to specification development has been one of “write once, adopt everywhere”,
390 with the deliverables being a set of specifications that carry the IEEE POSIX designation, The
391 Open Group’s Technical Standard designation, and an ISO/IEC designation.

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

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

400 Background

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

404 Conceptually, this standard describes a set of fundamental services needed for the efficient
405 construction of application programs. Access to these services has been provided by defining an
406 interface, using the C programming language, a command interpreter, and common utility
407 programs that establish standard semantics and syntax. Since this interface enables application
408 writers to write portable applications—it was developed with that goal in mind—it has been
409 designated POSIX,² an acronym for Portable Operating System Interface.

410 Although originated to refer to the original IEEE Std 1003.1-1988, the name POSIX more
411 correctly refers to a *family* of related standards: IEEE Std 1003.*n* and the parts of ISO/IEC 9945.
412 In earlier editions of the IEEE standard, the term POSIX was used as a synonym for
413 IEEE Std 1003.1-1988. A preferred term, POSIX.1, emerged. This maintained the advantages of

-
- 414 1. The Austin Group is named after the location of the inaugural meeting held at the IBM facility in Austin, Texas in September 1998.
415 2. The name POSIX was suggested by Richard Stallman. It is expected to be pronounced *pahz-icks*, as in *positive*, not *poh-six*, or other
416 variations. The pronunciation has been published in an attempt to promulgate a standardized way of referring to a standard operating
417 system interface.

Introduction

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.

- 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

460 have been added as optional extensions.

- 461 • Broadly Implementable

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

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

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

- 473 • Minimal Changes to Historical Implementations

474 When the original version—IEEE Std 1003.1-1988—was published, there were no known
475 historical implementations that did not have to change. However, there was a broad
476 consensus on a set of functions, types, definitions, and concepts that formed an interface
477 that was common to most historical implementations.

478 The adoption of the 1988 and 1990 IEEE system interface standards, the 1992 IEEE shell
479 and utilities standard, the various Open Group (formerly X/Open) specifications, and the
480 2001 Edition of this standard and its technical corrigenda have consolidated this
481 consensus, and this revision reflects the significantly increased level of consensus arrived
482 at since the original versions. The authors of the original versions tried, as much as
483 possible, to follow the principles below when creating new specifications:

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

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

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

- 499 • Minimal Changes to Existing Application Code

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

*Introduction*503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540**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.

XSI option

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

541

Participants

542

543

544

IEEE Std 1003.1-200x was prepared by the Austin Group, sponsored by the Portable Applications Standards Committee of the IEEE Computer Society, The Open Group, and ISO/SC22.

545

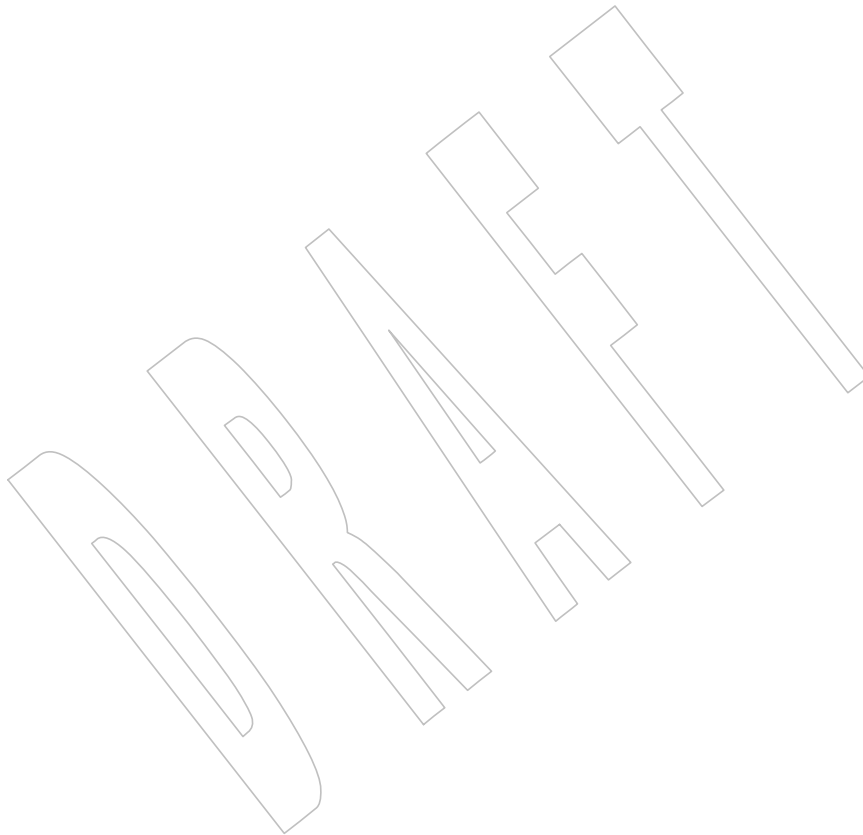
546

547

Notes to Reviewers

This section with side shading will not appear in the final copy. - Ed.

This section will be completed once the standard is approved.



Trademarks

548

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

553 1003.1™ is a trademark of the Institute of Electrical and Electronic Engineers, Inc.

554 AIX® is a registered trademark of IBM Corporation.

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

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

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

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

559 Boundaryless Information Flow™ and TOGAF™ are trademarks and Motif®, Making Standards
560 Work®, OSF/1®, The Open Group®, UNIX®, and the “X” device are registered trademarks of
561 The Open Group in the United States and other countries.

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

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

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

566

Acknowledgements

567

568

The contributions of the following organizations to the development of IEEE Std 1003.1-200x are gratefully acknowledged:

569

570

- 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.

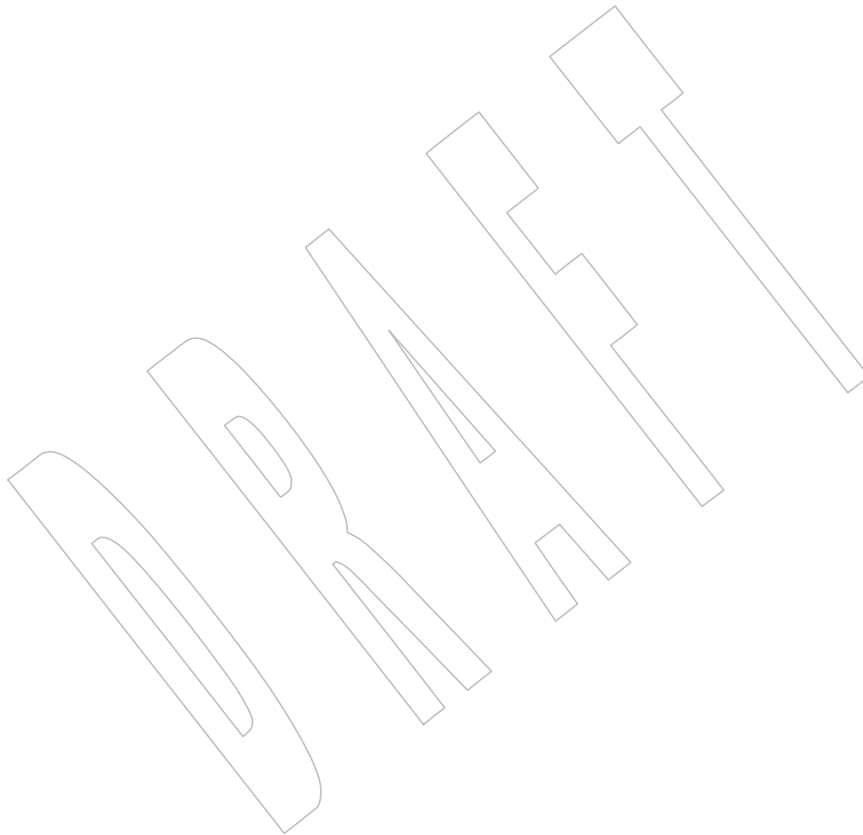
571

- ISO/IEC JTC 1/SC 22/WG 14 C Language Committee

572

573

This standard was prepared by the Austin Group, a joint working group of the IEEE, The Open Group, and ISO/IEC JTC 1/SC 22.



Referenced Documents

574

575

Normative References

576

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

577

Informative References

578

The following documents are referenced in this standard:

579

1984 /usr/group Standard

580

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

581

Almasi and Gottlieb

582

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

583

584

ANSI C

585

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

586

587

ANSI X3.226-1994

588

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

589

590

Brawer

591

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

592

593

DeRemer and Pennello Article

594

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

595

596

Draft ANSI X3J11.1

597

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

598

FIPS 151-1

599

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

600

601

FIPS 151-2

602

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

603

604

HP-UX Manual

605

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

606

IEC 60559: 1989

607

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

608

609

IEEE Std 754-1985

610

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

- 611 IEEE Std 854-1987
612 IEEE Std 854-1987, IEEE Standard for Radix-Independent Floating-Point Arithmetic.
- 613 IEEE Std 1003.9-1992
614 IEEE Std 1003.9-1992, IEEE Standard for Information Technology — POSIX FORTRAN 77
615 Language Interfaces — Part 1: Binding for System Application Program Interface API.
- 616 IETF RFC 791
617 Internet Protocol, Version 4 (IPv4), September 1981.
- 618 IETF RFC 819
619 The Domain Naming Convention for Internet User Applications, Z. Su, J. Postel, August
620 1982.
- 621 IETF RFC 822
622 Standard for the Format of ARPA Internet Text Messages, D.H. Crocker, August 1982.
- 623 IETF RFC 919
624 Broadcasting Internet Datagrams, J. Mogul, October 1984.
- 625 IETF RFC 920
626 Domain Requirements, J. Postel, J. Reynolds, October 1984.
- 627 IETF RFC 921
628 Domain Name System Implementation Schedule, J. Postel, October 1984.
- 629 IETF RFC 922
630 Broadcasting Internet Datagrams in the Presence of Subnets, J. Mogul, October 1984.
- 631 IETF RFC 1034
632 Domain Names — Concepts and Facilities, P. Mockapetris, November 1987.
- 633 IETF RFC 1035
634 Domain Names — Implementation and Specification, P. Mockapetris, November 1987.
- 635 IETF RFC 1123
636 Requirements for Internet Hosts — Application and Support, R. Braden, October 1989.
- 637 IETF RFC 1886
638 DNS Extensions to Support Internet Protocol, Version 6 (IPv6), C. Huitema, S. Thomson,
639 December 1995.
- 640 IETF RFC 2045
641 Multipurpose Internet Mail Extensions (MIME), Part 1: Format of Internet Message Bodies,
642 N. Freed, N. Borenstein, November 1996.
- 643 IETF RFC 2181
644 Clarifications to the DNS Specification, R. Elz, R. Bush, July 1997.
- 645 IETF RFC 2373
646 Internet Protocol, Version 6 (IPv6) Addressing Architecture, S. Deering, R. Hinden, July
647 1998.
- 648 IETF RFC 2460
649 Internet Protocol, Version 6 (IPv6), S. Deering, R. Hinden, December 1998.
- 650 Internationalisation Guide
651 Guide, July 1993, Internationalisation Guide, Version 2 (ISBN: 1-859120-02-4, G304),
652 published by The Open Group.

Referenced Documents

- 653 ISO C (1990)
 654 ISO/IEC 9899:1990, Programming Languages — C, including Amendment 1:1995 (E), C
 655 Integrity (Multibyte Support Extensions (MSE) for ISO C).
- 656 ISO 2375:1985
 657 ISO 2375:1985, Data Processing — Procedure for Registration of Escape Sequences.
- 658 ISO 8652:1987
 659 ISO 8652:1987, Programming Languages — Ada (technically identical to ANSI standard
 660 1815A-1983).
- 661 ISO/IEC 1539:1990
 662 ISO/IEC 1539:1990, Information Technology — Programming Languages — Fortran
 663 (technically identical to the ANSI X3.9-1978 standard [FORTRAN 77]).
- 664 ISO/IEC 4873:1991
 665 ISO/IEC 4873:1991, Information Technology — ISO 8-bit Code for Information Interchange
 666 — Structure and Rules for Implementation.
- 667 ISO/IEC 6429:1992
 668 ISO/IEC 6429:1992, Information Technology — Control Functions for Coded Character
 669 Sets.
- 670 ISO/IEC 6937:1994
 671 ISO/IEC 6937:1994, Information Technology — Coded Character Set for Text
 672 Communication — Latin Alphabet.
- 673 ISO/IEC 8802-3:1996
 674 ISO/IEC 8802-3:1996, Information Technology — Telecommunications and Information
 675 Exchange Between Systems — Local and Metropolitan Area Networks — Specific
 676 Requirements — Part 3: Carrier Sense Multiple Access with Collision Detection
 677 (CSMA/CD) Access Method and Physical Layer Specifications.
- 678 ISO/IEC 8859
 679 ISO/IEC 8859, Information Technology — 8-Bit Single-Byte Coded Graphic Character Sets:
 680 Part 1: Latin Alphabet No. 1
 681 Part 2: Latin Alphabet No. 2
 682 Part 3: Latin Alphabet No. 3
 683 Part 4: Latin Alphabet No. 4
 684 Part 5: Latin/Cyrillic Alphabet
 685 Part 6: Latin/Arabic Alphabet
 686 Part 7: Latin/Greek Alphabet
 687 Part 8: Latin/Hebrew Alphabet
 688 Part 9: Latin Alphabet No. 5
 689 Part 10: Latin Alphabet No. 6
 690 Part 11: Latin/Thai Alphabet
 691 Part 13: Latin Alphabet No. 7
 692 Part 14: Latin Alphabet No. 8
 693 Part 15: Latin Alphabet No. 9
 694 Part 16: Latin Alphabet No. 10
- 695 ISO POSIX-1:1996
 696 ISO/IEC 9945-1:1996, Information Technology — Portable Operating System Interface
 697 (POSIX) — Part 1: System Application Program Interface (API) [C Language] (identical to
 698 ANSI/IEEE Std 1003.1-1996). Incorporating ANSI/IEEE Stds 1003.1-1990, 1003.1b-1993,
 699 1003.1c-1995, and 1003.1i-1995.

- 700 ISO POSIX-2: 1993
 701 ISO/IEC 9945-2: 1993, Information Technology — Portable Operating System Interface
 702 (POSIX) — Part 2: Shell and Utilities (identical to ANSI/IEEE Std 1003.2-1992, as amended
 703 by ANSI/IEEE Std 1003.2a-1992).
- 704 Issue 1
 705 X/Open Portability Guide, July 1985 (ISBN: 0-444-87839-4).
- 706 Issue 2
 707 X/Open Portability Guide, January 1987:
- 708 • Volume 1: XVS Commands and Utilities (ISBN: 0-444-70174-5)
 - 709 • Volume 2: XVS System Calls and Libraries (ISBN: 0-444-70175-3)
- 710 Issue 3
 711 X/Open Specification, 1988, 1989, February 1992:
- 712 • Commands and Utilities, Issue 3 (ISBN: 1-872630-36-7, C211); this specification was
 713 formerly X/Open Portability Guide, Issue 3, Volume 1, January 1989, XSI Commands
 714 and Utilities (ISBN: 0-13-685835-X, XO/XPG/89/002)
 - 715 • System Interfaces and Headers, Issue 3 (ISBN: 1-872630-37-5, C212); this specification
 716 was formerly X/Open Portability Guide, Issue 3, Volume 2, January 1989, XSI System
 717 Interface and Headers (ISBN: 0-13-685843-0, XO/XPG/89/003)
 - 718 • Curses Interface, Issue 3, contained in Supplementary Definitions, Issue 3
 719 (ISBN: 1-872630-38-3, C213), Chapters 9 to 14 inclusive; this specification was formerly
 720 X/Open Portability Guide, Issue 3, Volume 3, January 1989, XSI Supplementary
 721 Definitions (ISBN: 0-13-685850-3, XO/XPG/89/004)
 - 722 • Headers Interface, Issue 3, contained in Supplementary Definitions, Issue 3
 723 (ISBN: 1-872630-38-3, C213), Chapter 19, Cpio and Tar Headers; this specification was
 724 formerly X/Open Portability Guide Issue 3, Volume 3, January 1989, XSI
 725 Supplementary Definitions (ISBN: 0-13-685850-3, XO/XPG/89/004)
- 726 Issue 4
 727 CAE Specification, July 1992, published by The Open Group:
- 728 • System Interface Definitions (XBD), Issue 4 (ISBN: 1-872630-46-4, C204)
 - 729 • Commands and Utilities (XCU), Issue 4 (ISBN: 1-872630-48-0, C203)
 - 730 • System Interfaces and Headers (XSH), Issue 4 (ISBN: 1-872630-47-2, C202)
- 731 Issue 4, Version 2
 732 CAE Specification, August 1994, published by The Open Group:
- 733 • System Interface Definitions (XBD), Issue 4, Version 2 (ISBN: 1-85912-036-9, C434)
 - 734 • Commands and Utilities (XCU), Issue 4, Version 2 (ISBN: 1-85912-034-2, C436)
 - 735 • System Interfaces and Headers (XSH), Issue 4, Version 2 (ISBN: 1-85912-037-7, C435)
- 736 Issue 5
 737 Technical Standard, February 1997, published by The Open Group:
- 738 • System Interface Definitions (XBD), Issue 5 (ISBN: 1-85912-186-1, C605)
 - 739 • Commands and Utilities (XCU), Issue 5 (ISBN: 1-85912-191-8, C604)

Referenced Documents

- 740 • System Interfaces and Headers (XSH), Issue 5 (ISBN: 1-85912-181-0, C606)
- 741 Issue 6
- 742 Technical Standard, April 2004, published by The Open Group:
- 743 • Base Definitions (XBD), Issue 6 (ISBN: 1-931624-43-7, C046)
- 744 • System Interfaces (XSH), Issue 6 (ISBN: 1-931624-44-5, C047)
- 745 • Shell and Utilities (XCU), Issue 6 (ISBN: 1-931624-45-3, C048)
- 746 Knuth Article
- 747 Knuth, Donald E., *On the Translation of Languages from Left to Right*, Information and Control,
748 Volume 8, No. 6, October 1965.
- 749 KornShell
- 750 Bolsky, Morris I. and Korn, David G., *The New KornShell Command and Programming
751 Language*, March 1995, Prentice Hall.
- 752 MSE Working Draft
- 753 Working draft of ISO/IEC 9899:1990/Add3:Draft, Addendum 3 — Multibyte Support
754 Extensions (MSE) as documented in the ISO Working Paper SC22/WG14/N205 dated 31
755 March 1992.
- 756 POSIX.0: 1995
- 757 IEEE Std 1003.0-1995, IEEE Guide to the POSIX Open System Environment (OSE) (identical
758 to ISO/IEC TR 14252).
- 759 POSIX.1: 1988
- 760 IEEE Std 1003.1-1988, IEEE Standard for Information Technology — Portable Operating
761 System Interface (POSIX) — Part 1: System Application Program Interface (API) [C
762 Language].
- 763 POSIX.1: 1990
- 764 IEEE Std 1003.1-1990, IEEE Standard for Information Technology — Portable Operating
765 System Interface (POSIX) — Part 1: System Application Program Interface (API) [C
766 Language].
- 767 POSIX.1a
- 768 P1003.1a, Standard for Information Technology — Portable Operating System Interface
769 (POSIX) — Part 1: System Application Program Interface (API) — (C Language)
770 Amendment.
- 771 POSIX.1d: 1999
- 772 IEEE Std 1003.1d-1999, IEEE Standard for Information Technology — Portable Operating
773 System Interface (POSIX) — Part 1: System Application Program Interface (API) —
774 Amendment 4: Additional Realtime Extensions [C Language].
- 775 POSIX.1g: 2000
- 776 IEEE Std 1003.1g-2000, IEEE Standard for Information Technology — Portable Operating
777 System Interface (POSIX) — Part 1: System Application Program Interface (API) —
778 Amendment 6: Protocol-Independent Interfaces (PII).
- 779 POSIX.1j: 2000
- 780 IEEE Std 1003.1j-2000, IEEE Standard for Information Technology — Portable Operating
781 System Interface (POSIX) — Part 1: System Application Program Interface (API) —
782 Amendment 5: Advanced Realtime Extensions [C Language].

- 783 POSIX.1q: 2000
784 IEEE Std 1003.1q-2000, IEEE Standard for Information Technology — Portable Operating
785 System Interface (POSIX) — Part 1: System Application Program Interface (API) —
786 Amendment 7: Tracing [C Language].
- 787 POSIX.2b
788 P1003.2b, Standard for Information Technology — Portable Operating System Interface
789 (POSIX) — Part 2: Shell and Utilities — Amendment.
- 790 POSIX.2d:-1994
791 IEEE Std 1003.2d-1994, IEEE Standard for Information Technology — Portable Operating
792 System Interface (POSIX) — Part 2: Shell and Utilities — Amendment 1: Batch Environment.
- 793 POSIX.13:-1998
794 IEEE Std 1003.13:1998, IEEE Standard for Information Technology — Standardized
795 Application Environment Profile (AEP) — POSIX Realtime Application Support.
- 796 Sarwate Article
797 Sarwate, Dilip V., *Computation of Cyclic Redundancy Checks via Table Lookup*, Communications
798 of the ACM, Volume 30, No. 8, August 1988.
- 799 Sprunt, Sha, and Lehoczky
800 Sprunt, B., Sha, L., and Lehoczky, J.P., *Aperiodic Task Scheduling for Hard Real-Time Systems*,
801 The Journal of Real-Time Systems, Volume 1, 1989, Pages 27-60.
- 802 SVID, Issue 1
803 American Telephone and Telegraph Company, System V Interface Definition (SVID), Issue
804 1; Morristown, NJ, UNIX Press, 1985.
- 805 SVID, Issue 2
806 American Telephone and Telegraph Company, System V Interface Definition (SVID), Issue
807 2; Morristown, NJ, UNIX Press, 1986.
- 808 SVID, Issue 3
809 American Telephone and Telegraph Company, System V Interface Definition (SVID), Issue
810 3; Morristown, NJ, UNIX Press, 1989.
- 811 The AWK Programming Language
812 Aho, Alfred V., Kernighan, Brian W., and Weinberger, Peter J., *The AWK Programming
813 Language*, Reading, MA, Addison-Wesley 1988.
- 814 UNIX Programmer's Manual
815 American Telephone and Telegraph Company, *UNIX Time-Sharing System: UNIX
816 Programmer's Manual*, 7th Edition, Murray Hill, NJ, Bell Telephone Laboratories, January
817 1979.
- 818 XNS, Issue 4
819 CAE Specification, August 1994, Networking Services, Issue 4 (ISBN: 1-85912-049-0, C438),
820 published by The Open Group.
- 821 XNS, Issue 5
822 CAE Specification, February 1997, Networking Services, Issue 5 (ISBN: 1-85912-165-9, C523),
823 published by The Open Group.
- 824 XNS, Issue 5.2
825 Technical Standard, January 2000, Networking Services (XNS), Issue 5.2
826 (ISBN: 1-85912-241-8, C808), published by The Open Group.

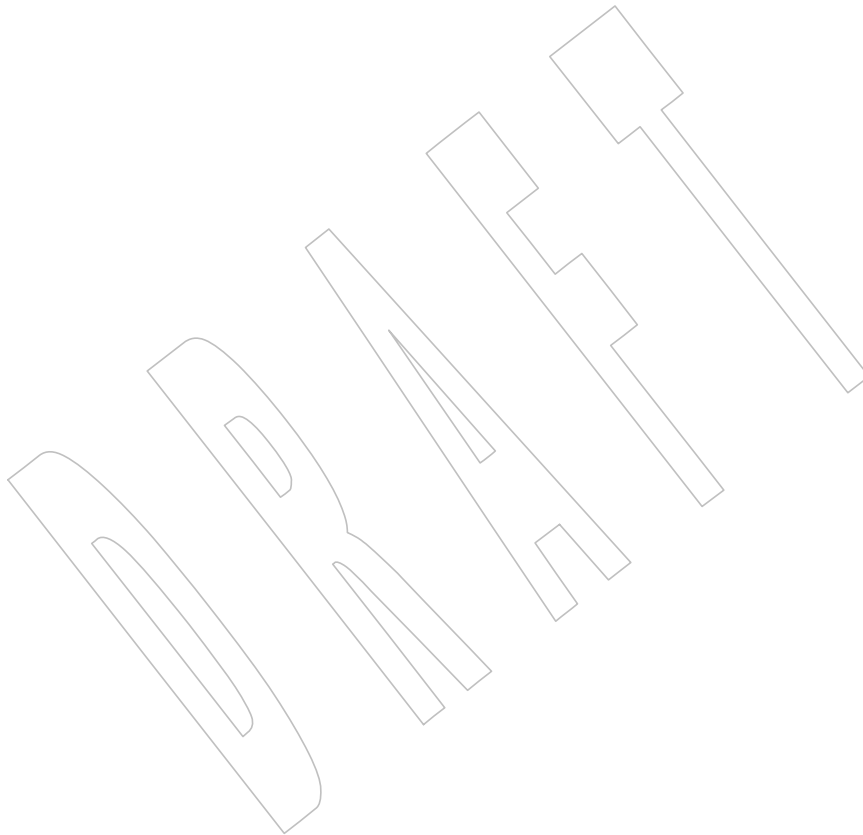
Referenced Documents

- 827 X/Open Curses, Issue 4, Version 2
828 CAE Specification, May 1996, X/Open Curses, Issue 4, Version 2 (ISBN: 1-85912-171-3,
829 C610), published by The Open Group.
- 830 Yacc
831 *Yacc: Yet Another Compiler Compiler*, Stephen C. Johnson, 1978.

832 Source Documents

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

- 834 AIX 3.2 Manual
835 AIX Version 3.2 For RISC System/6000, Technical Reference: Base Operating System and
836 Extensions, 1990, 1992 (Part No. SC23-2382-00).
- 837 OSF/1
838 OSF/1 Programmer's Reference, Release 1.2 (ISBN: 0-13-020579-6).
- 839 OSF AES
840 Application Environment Specification (AES) Operating System Programming Interfaces
841 Volume, Revision A (ISBN: 0-13-043522-8).
- 842 System V Release 2.0
843 — UNIX System V Release 2.0 Programmer's Reference Manual (April 1984 - Issue 2).
844 — UNIX System V Release 2.0 Programming Guide (April 1984 - Issue 2).
- 845 System V Release 4.2
846 Operating System API Reference, UNIX[®] SVR4.2 (1992) (ISBN: 0-13-017658-3).



Rationale (Informative)

1

Part A:

2

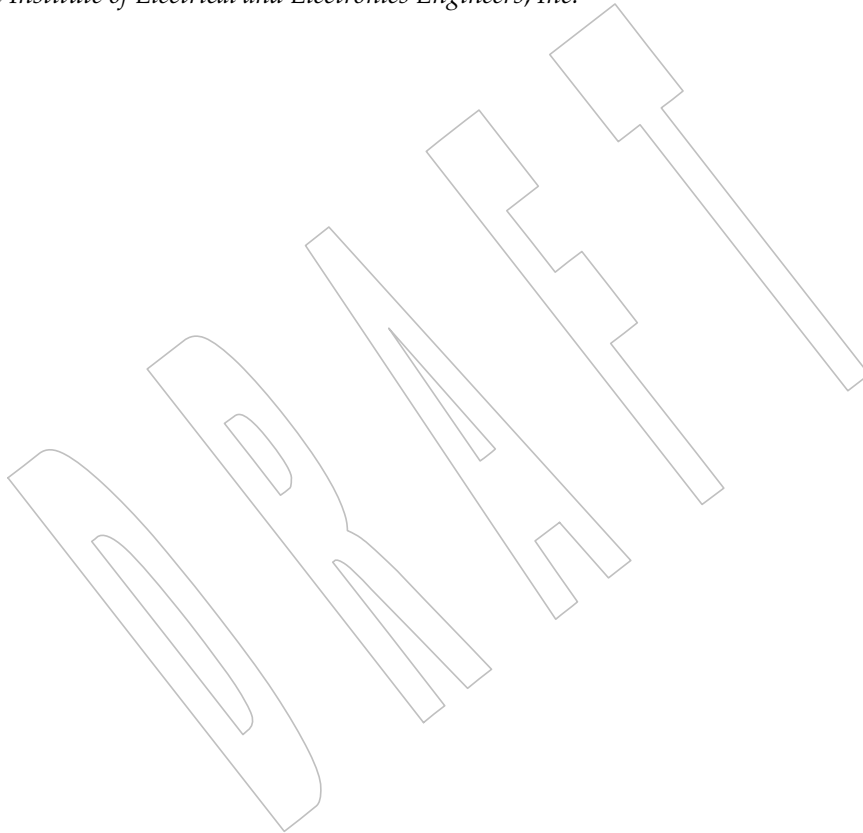
Base Definitions

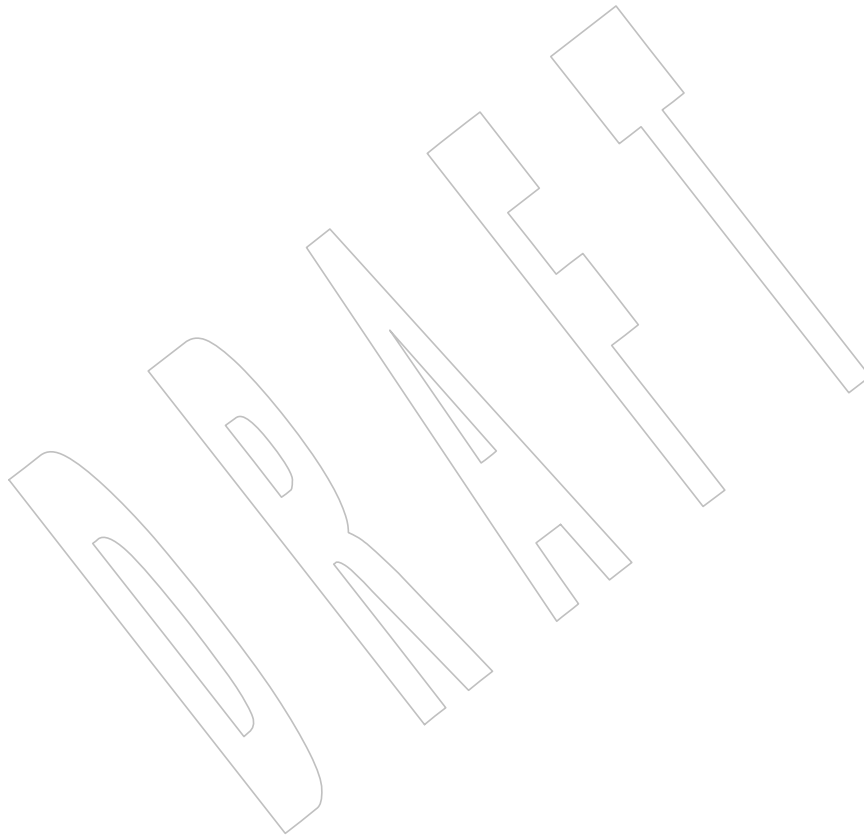
3

The Open Group

4

The Institute of Electrical and Electronics Engineers, Inc.





Rationale for Base Definitions

A.1 Introduction

A.1.1 Scope

IEEE Std 1003.1-200x is one of a family of standards known as POSIX. The family of standards extends to many topics; IEEE Std 1003.1-200x is known as POSIX.1 and consists of both operating system interfaces and shell and utilities. IEEE Std 1003.1-200x is technically identical to The Open Group Base Specifications, Issue 7.

Scope of IEEE Std 1003.1-200x

The (paraphrased) goals of this development were to revise the single document that is ISO/IEC 9945:2003 Parts 1 through 4, IEEE Std 1003.1, 2004 Edition, and the appropriate parts of The Open Group Single UNIX Specification, Version 3. This work has been undertaken by the Austin Group, a joint working group of IEEE, The Open Group, and ISO/IEC JTC 1/SC 22.

The following are the base documents in this revision:

- IEEE Std 1003.1, 2004 Edition
- ISO/IEC 9899:1999, Programming Languages — C (including ISO/IEC 9899:1999/Cor.1:2001(E) and ISO/IEC 9899:1999/Cor.2:2004(E))
- The Open Group Extended API Sets, Parts 1 through 4

The revision has addressed the following areas:

- Issues raised by Austin Group defect reports, IEEE Interpretations against IEEE Std 1003.1, and ISO/IEC defect reports against ISO/IEC 9945

The repository of interpretations can be accessed at www.opengroup.org/austin/interps.

- Issues raised in corrigenda for The Open Group Technical Standards and working group resolutions from The Open Group
- Issues arising from ISO TR 24715: 2006, Conflicts between POSIX and the LSB

This is a Type 3 informative technical report highlighting differences between the LSB 3.1 and the 2004 edition of this standard.

- Changes to make the text self-consistent with the additional material merged

The new material merged has come from the The Open Group Extended API Sets Parts 1 through 4. A list of the new interfaces is included in [Section B.1.4](#) (on page 79).

- Features, marked legacy or obsolescent in the base documents, have been considered for withdrawal in the revision

See [Section B.1.4](#) and [Section C.1.4](#) (on page 219).

- A review and reorganization of the options within the standard

This has included marking the following options obsolescent:

- Batch Environment Services and Utilities
- Tracing
- XSI STREAMS

The UUCP Utilities option is a new option for this revision.

Functionality from the following former options is now mandatory in this version:

AIO	_POSIX_ASYNCHRONOUS_IO (Asynchronous Input and Output)
BAR	_POSIX_BARRIERS (Barriers)
CS	_POSIX_CLOCK_SELECTION (Clock Selection)
MF	_POSIX_MAPPED_FILES (Memory Mapped Files)
MPR	_POSIX_MEMORY_PROTECTION (Memory Protection)
RTS	_POSIX_REALTIME_SIGNALS (Realtime Signals Extension)
RWL	_POSIX_READER_WRITER_LOCKS (Read-Write Locks)
SEM	_POSIX_SEMAPHORES (Semaphores)
SPI	_POSIX_SPIN_LOCKS (Spin Locks)
THR	_POSIX_THREADS (Threads)
TMO	_POSIX_TIMEOUTS (Timeouts)
TMR	_POSIX_TIMERS (Timers)
TSF	_POSIX_THREAD_SAFE_FUNCTIONS (Thread-Safe Functions)

- Alignment with the ISO/IEC 9899:1999 standard, including Technical Corrigendum 2
- A review of the use of fixed path filenames within the standard

For example, the *at*, *batch*, and *crontab* utilities previously had a requirement for use of the directory **/usr/lib/cron**.

The following were requirements on IEEE Std 1003.1-200x:

- Backward-compatibility

For interfaces carried forward, it was agreed that there should be no breakage of functionality in the existing base documents. All strictly conforming applications will be conforming but not necessarily strictly conforming to the revised standard. The goal is for system implementations to be able to support the existing and revised standards simultaneously.

- Architecture and *n*-bit-neutral

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

- Extensibility

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

77 POSIX.1 and the ISO C Standard

78 The standard developers believed it essential for a programmer to have a single complete
79 reference place, but recognized that deference to the formal standard has to be addressed for the
80 duplicate interface definitions between the ISO C standard and IEEE Std 1003.1-200x.

81 Where an interface has a version in the ISO C standard, the DESCRIPTION section describes the
82 relationship to the ISO C standard and markings are included as appropriate to show where the
83 ISO C standard has been extended in the text.

84 A block of text is included at the start of each affected reference page stating whether the page is
85 aligned with the ISO C standard or extended. Each page has been parsed for additions beyond
86 the ISO C standard (that is, including both POSIX and UNIX extensions), and these extensions
87 are marked as CX extensions (for C extensions).

88 FIPS Requirements

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

92 The following restrictions are carried forward from IEEE Std 1003.1-2001 in order to align with
93 FIPS 151-2 requirements:

- 94 • The implementation supports `_POSIX_CHOWN_RESTRICTED`.
- 95 • The limit `{NGROUPS_MAX}` is greater than or equal to 8.
- 96 • The implementation supports the setting of the group ID of a file (when it is created) to
97 that of the parent directory.
- 98 • The implementation supports `_POSIX_SAVED_IDS`.
- 99 • The implementation supports `_POSIX_VDISABLE`.
- 100 • The implementation supports `_POSIX_JOB_CONTROL`.
- 101 • The implementation supports `_POSIX_NO_TRUNC`.
- 102 • The `read()` function returns the number of bytes read when interrupted by a signal and
103 does not return `-1`.
- 104 • The `write()` function returns the number of bytes written when interrupted by a signal and
105 does not return `-1`.
- 106 • In the environment for the login shell, the environment variables `LOGNAME` and `HOME`
107 are defined and have the properties described in IEEE Std 1003.1-200x.
- 108 • The value of `{CHILD_MAX}` is greater than or equal to 25.
- 109 • The value of `{OPEN_MAX}` is greater than or equal to 20.
- 110 • The implementation supports the functionality associated with the symbols `CS7`, `CS8`,
111 `CSTOPB`, `PARODD`, and `PARENB` defined in `<termios.h>`.

112 A.1.2 Conformance

113 See [Section A.2](#) (on page 9).

114 A.1.3 Normative References

115 There is no additional rationale provided for this section.

116 A.1.4 Terminology

117 The meanings specified in IEEE Std 1003.1-200x for the words *shall*, *should*, and *may* are
118 mandated by ISO/IEC directives.

119 In the Rationale (Informative) volume of IEEE Std 1003.1-200x, the words *shall*, *should*, and *may*
120 are sometimes used to illustrate similar usages in IEEE Std 1003.1-200x. However, the rationale
121 itself does not specify anything regarding implementations or applications.

122 conformance document

123 As a practical matter, the conformance document is effectively part of the system
124 documentation. Conformance documents are distinguished by IEEE Std 1003.1-200x so that
125 they can be referred to distinctly.

126 implementation-defined

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

130 may

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

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

136 shall

137 Declarative sentences are sometimes used in IEEE Std 1003.1-200x as if they included the
138 word *shall*, and facilities thus specified are no less required. For example, the two
139 statements:

- 140 1. The *foo()* function shall return zero.
- 141 2. The *foo()* function returns zero.

142 are meant to be exactly equivalent.

143 should

144 In IEEE Std 1003.1-200x, the word *should* does not usually apply to the implementation, but
145 rather to the application. Thus, the important words regarding implementations are *shall*,
146 which indicates requirements, and *may*, which indicates options.

147 obsolescent

148 The term “obsolescent” means “do not use this feature in new applications”. A feature
149 noted as obsolescent is supported by all implementations, but may be removed in a future
150 version; new applications should not use these features. The obsolescence concept is not an
151 ideal solution, but was used as a method of increasing consensus: many more objections
152 would be heard from the user community if some of these historical features were suddenly
153 withdrawn without the grace period obsolescence implies. The phrase “may be removed in
154 a future version” implies that the result of that consideration might in fact keep those
155 features indefinitely if the predominance of applications do not migrate away from them

156 quickly.

157 **legacy**

158 The term “legacy” was included in previous versions of this standard but is no longer used
159 in the current version.

160 **system documentation**

161 The system documentation should normally describe the whole of the implementation,
162 including any extensions provided by the implementation. Such documents normally
163 contain information at least as detailed as the specifications in IEEE Std 1003.1-200x. Few
164 requirements are made on the system documentation, but the term is needed to avoid a
165 dangling pointer where the conformance document is permitted to point to the system
166 documentation.

167 **undefined**

168 See *implementation-defined*.

169 **unspecified**

170 See *implementation-defined*.

171 The definitions for “unspecified” and “undefined” appear nearly identical at first
172 examination, but are not. The term “unspecified” means that a conforming application may
173 deal with the unspecified behavior, and it should not care what the outcome is. The term
174 “undefined” says that a conforming application should not do it because no definition is
175 provided for what it does (and implicitly it would care what the outcome was if it tried it).
176 It is important to remember, however, that if the syntax permits the statement at all, it must
177 have some outcome in a real implementation.

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

182 This also implies that an implementation, or another standard, could specify or define the
183 result in a useful fashion. The terms apply to IEEE Std 1003.1-200x specifically.

184 The term “implementation-defined” implies requirements for documentation that are not
185 required for “undefined” (or “unspecified”). Where there is no need for a conforming
186 program to know the definition, the term “undefined” is used, even though
187 “implementation-defined” could also have been used in this context. There could be a
188 fourth term, specifying “this standard does not say what this does; it is acceptable to define
189 it in an implementation, but it does not need to be documented”, and undefined would
190 then be used very rarely for the few things for which any definition is not useful. In
191 particular, implementation-defined is used where it is believed that certain classes of
192 application will need to know such details to determine whether the application can be
193 successfully ported to the implementation. Such applications are not always strictly
194 portable, but nevertheless are common and useful; often the requirements met by the
195 application cannot be met without dealing with the issues implied by “implementation-
196 defined”. In some places the text refers to facilities supplied by the implementation that are
197 outside the standard as implementation-supplied or implementation-provided. This is not
198 intended to imply a requirement for documentation. If it were, the term “implementation-
199 defined” would have been used.

200 In many places IEEE Std 1003.1-200x is silent about the behavior of some possible construct.
201 For example, a variable may be defined for a specified range of values and behaviors are
202 described for those values; nothing is said about what happens if the variable has any other
203 value. That kind of silence can imply an error in the standard, but it may also imply that the

204 standard was intentionally silent and that any behavior is permitted. There is a natural
 205 tendency to infer that if the standard is silent, a behavior is prohibited. That is not the intent.
 206 Silence is intended to be equivalent to the term “unspecified”.

207 The term “application” is not defined in IEEE Std 1003.1-200x; it is assumed to be a part of
 208 general computer science terminology.

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

211 **macro**

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

215 **symbolic name**

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

219 **symbolic constant**

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

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

223 **A.1.5 Portability**

224 To aid the identification of options within IEEE Std 1003.1-200x, a notation consisting of margin
 225 codes and shading is used. This is based on the notation used in previous revisions of The Open
 226 Group Base specifications.

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

232 **A.1.5.1 Codes**

233 This section includes codes for options defined in the Base Definitions volume of
 234 IEEE Std 1003.1-200x, Section 2.1.6, Options, and the following additional codes for other
 235 purposes:

236 **CX** This margin code is used to denote extensions beyond the ISO C standard. For
 237 interfaces that are duplicated between IEEE Std 1003.1-200x and the ISO C standard, a
 238 CX introduction block describes the nature of the duplication, with any extensions
 239 appropriately CX marked and shaded.

240 Where an interface is added to an ISO C standard header, within the header the
 241 interface has an appropriate margin marker and shading (for example, CX, XSI, TSE,
 242 and so on) and the same marking appears on the reference page in the SYNOPSIS
 243 section. This enables a programmer to easily identify that the interface is extending an
 244 ISO C standard header.

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

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

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

252 A.1.5.2 Margin Code Notation

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

258 A.2 Conformance

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

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

262 A.2.1 Implementation Conformance

263 These definitions allow application developers to know what to depend on in an
 264 implementation.

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

269 A.2.1.1 Requirements

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

273 A.2.1.2 Documentation

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

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

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

288 Also, interrelated file activities, where multiple files and/or directories are updated, or where
 289 space is allocated or released in the file system structures, can leave inconsistencies in the
 290 relationship between data in the various files and directories, or in the file system itself. Such

291 inconsistencies can break applications that expect updates to occur in a specific sequence, so that
 292 updates in one place correspond with related updates in another place.

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

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

306 A.2.1.3 POSIX Conformance

307 This really means conformance to the base standard; however, since this document includes the
 308 core material of the Single UNIX Specification, the standard developers decided that it was
 309 appropriate to segment the conformance requirements into two, the former for the base
 310 standard, and the latter for the Single UNIX Specification (denoted XSI Conformance).

311 Within POSIX.1 there are some symbolic constants that, if defined to a certain value or range of
 312 values, indicate that a certain option is enabled. Other symbolic constants exist in POSIX.1 for
 313 other reasons.

314 As part of the revision some features that were previously optional have been made mandatory.
 315 For backwards compatibility, the symbolic constants associated with the option are still required
 316 now with fixed allowable ranges or values. The following options from the previous version of
 317 this standard are now mandatory:

318 `_POSIX_ASYNCHRONOUS_IO`
 319 `_POSIX_BARRIERS`
 320 `_POSIX_CLOCK_SELECTION`
 321 `_POSIX_MAPPED_FILES`
 322 `_POSIX_MEMORY_PROTECTION`
 323 `_POSIX_READER_WRITER_LOCKS`
 324 `_POSIX_REALTIME_SIGNALS`
 325 `_POSIX_SEMAPHORES`
 326 `_POSIX_SPIN_LOCKS`
 327 `_POSIX_THREAD_SAFE_FUNCTIONS`
 328 `_POSIX_THREADS`
 329 `_POSIX_TIMEOUTS`
 330 `_POSIX_TIMERS`

331 A POSIX-conformant system may support the XSI option required by the Single UNIX
 332 Specification. This was intentional since the standard developers intend them to be upwards-
 333 compatible, so that a system conforming to the Single UNIX Specification can also conform to
 334 the base standard at the same time.

335 A.2.1.4 XSI Conformance

336 This section is included to describe the conformance requirements for the base volumes of the
337 Single UNIX Specification.

338 XSI conformance can be thought of as a profile, selecting certain options from
339 IEEE Std 1003.1-200x.

340 A.2.1.5 Option Groups

341 The concept of “Option Groups” is included to allow collections of related functions or options
342 to be grouped together. This has been used as follows: the “XSI Option Groups” have been
343 created to allow super-options, collections of underlying options and related functions, to be
344 collectively supported by XSI-conforming systems.

345 The standard developers considered the matter of subprofiling and decided it was better to
346 include an enabling mechanism rather than detailed normative requirements. A set of
347 subprofiling options was developed and included later in this volume of IEEE Std 1003.1-200x as
348 an informative illustration.

349 **Subprofiling Considerations**

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

352 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)

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

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

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

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

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

373 Another example:

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

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

```

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

```

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

384 A.2.1.6 Options

385 The final subsections within *Implementation Conformance* list the core options within
 386 IEEE Std 1003.1-200x. This includes both options for the System Interfaces volume of
 387 IEEE Std 1003.1-200x and the Shell and Utilities volume of IEEE Std 1003.1-200x.

388 A.2.2 Application Conformance

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

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

398 A.2.2.1 Strictly Conforming POSIX Application

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

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

403 A.2.2.2 Conforming POSIX Application

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

405 A.2.2.3 Conforming POSIX Application Using Extensions

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

410 A.2.2.4 Strictly Conforming XSI Application

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

413 A.2.2.5 Conforming XSI Application Using Extensions

414 Such applications may use extensions beyond the facilities defined by IEEE Std 1003.1-200x
 415 including the XSI option, but need to document the additional requirements.

416 A.2.3 Language-Dependent Services for the C Programming Language

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

422 A.2.4 Other Language-Related Specifications

423 There is no additional rationale provided for this section.

424 A.3 Definitions

425 The definitions in this section are stated so that they can be used as exact substitutes for the
 426 terms in text. They should not contain requirements or cross-references to sections within
 427 IEEE Std 1003.1-200x; that is accomplished by using an informative note. In addition, the term
 428 should not be included in its own definition. Where requirements or descriptions need to be
 429 addressed but cannot be included in the definitions, due to not meeting the above criteria, these
 430 occur in the General Concepts chapter.

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

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

437 Some definitions must allow extension to cover terms or facilities that are not explicitly
 438 mentioned in IEEE Std 1003.1-200x. For example, the definition of “Extended Security Controls”
 439 permits implementations beyond those defined in IEEE Std 1003.1-200x.

440 Some terms in the following list of notes do not appear in IEEE Std 1003.1-200x; these are
 441 marked suffixed with an asterisk (*). Many of them have been specifically excluded from
 442 IEEE Std 1003.1-200x because they concern system administration, implementation, or other
 443 issues that are not specific to the programming interface. Those are marked with a reason, such
 444 as “implementation-defined”.

445 Appropriate Privileges

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

460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502

Byte

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

According to the ISO/IEC 9899:1999 standard:

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

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

Character

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

Child Process

IEEE Std 1003.1-2001/Cor 2-2004, item XBD/TC2/D6/3 is applied, adding the `vfork()` function to those listed.

Clock Tick

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

Command

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

The terms "tool" and "application" can be thought of as being synonymous with "utility" from

503 the perspective of the operating system kernel. Tools, applications, and utilities historically have
504 run, typically, in processes above the kernel level. Tools and utilities historically have been a part
505 of the operating system non-kernel code and have performed system-related functions, such as
506 listing directory contents, checking file systems, repairing file systems, or extracting system
507 status information. Applications have not generally been a part of the operating system, and
508 they perform non-system-related functions, such as word processing, architectural design,
509 mechanical design, workstation publishing, or financial analysis. Utilities have most frequently
510 been provided by the operating system distributor, applications by third-party software
511 distributors, or by the users themselves. Nevertheless, IEEE Std 1003.1-200x does not
512 differentiate between tools, utilities, and applications when it comes to receiving services from
513 the system, a shell, or the standard utilities. (For example, the *xargs* utility invokes another
514 utility; it would be of fairly limited usefulness if the users could not run their own applications
515 in place of the standard utilities.) Utilities are not applications in the sense that they are not
516 themselves subject to the restrictions of IEEE Std 1003.1-200x or any other standard—there is no
517 requirement for *grep*, *stty*, or any of the utilities defined here to be any of the classes of
518 conforming applications.

519 **Column Positions**

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

526 In IEEE Std 1003.1-200x the term “column positions” has been defined to mean character—not
527 byte—positions in input files (such as “column position 7 of the FORTRAN input”). Output files
528 describe the column position in terms of the display width of the narrowest printable character
529 in the character set, adjusted to fit the characteristics of the output device. It is very possible that
530 *n* column positions will not be able to hold *n* characters in some character sets, unless all of those
531 characters are of the narrowest width. It is assumed that the implementation is aware of the
532 width of the various characters, deriving this information from the value of *LC_CTYPE*, and
533 thus can determine how many column positions to allot for each character in those utilities
534 where it is important.

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

539 **Controlling Terminal**

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

543 **Device Number***544 The concept is handled in *stat()* as *ID of device*.545 **Direct I/O**546 Historically, direct I/O refers to the system bypassing intermediate buffering, but may be
547 extended to cover implementation-defined optimizations.548 **Directory**549 The format of the directory file is implementation-defined and differs radically between
550 System V and 4.3 BSD. However, routines (derived from 4.3 BSD) for accessing directories and
551 certain constraints on the format of the information returned by those routines are described in
552 the **<dirent.h>** header.553 **Directory Entry**554 Throughout IEEE Std 1003.1-200x, the term “link” is used (about the *link()* function, for
555 example) in describing the objects that point to files from directories.556 **Display**557 The Shell and Utilities volume of IEEE Std 1003.1-200x assigns precise requirements for the
558 terms “display” and “write”. Some historical systems have chosen to implement certain utilities
559 without using the traditional file descriptor model. For example, the *vi* editor might employ
560 direct screen memory updates on a personal computer, rather than a *write()* system call. An
561 instance of user prompting might appear in a dialog box, rather than with standard error. When
562 the Shell and Utilities volume of IEEE Std 1003.1-200x uses the term “display”, the method of
563 outputting to the terminal is unspecified; many historical implementations use *termcap* or
564 *terminfo*, but this is not a requirement. The term “write” is used when the Shell and Utilities
565 volume of IEEE Std 1003.1-200x mandates that a file descriptor be used and that the output can
566 be redirected. However, it is assumed that when the writing is directly to the terminal (it has not
567 been redirected elsewhere), there is no practical way for a user or test suite to determine whether
568 a file descriptor is being used. Therefore, the use of a file descriptor is mandated only for the
569 redirection case and the implementation is free to use any method when the output is not
570 redirected. The verb *write* is used almost exclusively, with the very few exceptions of those
571 utilities where output redirection need not be supported: *tabs*, *talk*, *tput*, and *vi*.572 **Dot**573 The symbolic name *dot* is carefully used in POSIX.1 to distinguish the working directory
574 filename from a period or a decimal point.575 **Dot-Dot**576 Historical implementations permit the use of these filenames without their special meanings.
577 Such use precludes any meaningful use of these filenames by a Conforming POSIX.1
578 Application. Therefore, such use is considered an extension, the use of which makes an
579 implementation non-conforming; see also [Section A.4.12](#) (on page 38).

580 Epoch

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

585 FIFO Special File

586 See [Pipe](#) (on page 24).

587 File

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

589 File Classes

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

593 Filename

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

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

603 File System

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

609 Graphic Character

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

612 Group Database

613 IEEE Std 1003.1-2001/Cor 2-2004, item XBD/TC2/D6/4 is applied, removing the words “of
614 implementation-defined format”. See [User Database](#) (on page 32).

615 **Group File***616 Implementation-defined; see [User Database](#) (on page 32).617 **Historical Implementations***618 This refers to previously existing implementations of programming interfaces and operating
619 systems that are related to the interface specified by POSIX.1.620 **Hosted Implementation***621 This refers to a POSIX.1 implementation that is accomplished through interfaces from the
622 POSIX.1 services to some alternate form of operating system kernel services. Note that the line
623 between a hosted implementation and a native implementation is blurred, since most
624 implementations will provide some services directly from the kernel and others through some
625 indirect path. (For example, *fopen()* might use *open()*; or *mkfifo()* might use *mknod()*.) There is
626 no necessary relationship between the type of implementation and its correctness, performance,
627 and/or reliability.628 **Implementation***629 This term is generally used instead of its synonym, “system”, to emphasize the consequences of
630 decisions to be made by system implementors. Perhaps if no options or extensions to POSIX.1
631 were allowed, this usage would not have occurred.632 The term “specific implementation” is sometimes used as a synonym for “implementation”.
633 This should not be interpreted too narrowly; both terms can represent a relatively broad group
634 of systems. For example, a hardware vendor could market a very wide selection of systems that
635 all used the same instruction set, with some systems desktop models and others large multi-user
636 minicomputers. This wide range would probably share a common POSIX.1 operating system,
637 allowing an application compiled for one to be used on any of the others; this is a [*specific*]
638 *implementation*. However, such a wide range of machines probably has some differences
639 between the models. Some may have different clock rates, different file systems, different
640 resource limits, different network connections, and so on, depending on their sizes or intended
641 usages. Even on two identical machines, the system administrators may configure them
642 differently. Each of these different systems is known by the term “a specific instance of a specific
643 implementation”. This term is only used in the portions of POSIX.1 dealing with runtime
644 queries: *sysconf()* and *pathconf()*.645 **Incomplete Pathname***

646 Absolute pathname has been adequately defined.

647 **Job Control**648 In order to understand the job control facilities in POSIX.1 it is useful to understand how they
649 are used by a job control-cognizant shell to create the user interface effect of job control.650 While the job control facilities supplied by POSIX.1 can, in theory, support different types of
651 interactive job control interfaces supplied by different types of shells, there was historically one
652 particular interface that was most common when the standard was originally developed
653 (provided by BSD C Shell).654 This discussion describes that interface as a means of illustrating how the POSIX.1 job control
655 facilities can be used.656 Job control allows users to selectively stop (suspend) the execution of processes and continue
657 (resume) their execution at a later point. The user typically employs this facility via the

658 interactive interface jointly supplied by the terminal I/O driver and a command interpreter
659 (shell).

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

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

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

675 *Implementing Job Control Shells*

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

756 *Implementing Job Control Applications*

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

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

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

782 *Implementing Job Control Systems*

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

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

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

800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833**Kernel***

See [System Call*](#) (on page 31).

Library Routine*

See [System Call*](#) (on page 31).

Logical Device*

Implementation-defined.

Map

The definition of map is included to clarify the usage of mapped pages in the description of the behavior of process memory locking.

Memory-Resident

The term “memory-resident” is historically understood to mean that the so-called resident pages are actually present in the physical memory of the computer system and are immune from swapping, paging, copy-on-write faults, and so on. This is the actual intent of IEEE Std 1003.1-200x in the process memory locking section for implementations where this is logical. But for some implementations—primarily mainframes—actually locking pages into primary storage is not advantageous to other system objectives, such as maximizing throughput. For such implementations, memory locking is a “hint” to the implementation that the application wishes to avoid situations that would cause long latencies in accessing memory. Furthermore, there are other implementation-defined issues with minimizing memory access latencies that “memory residency” does not address—such as MMU reload faults. The definition attempts to accommodate various implementations while allowing conforming applications to specify to the implementation that they want or need the best memory access times that the implementation can provide.

Memory Object*

The term “memory object” usually implies shared memory. If the object is the same as a filename in the file system name space of the implementation, it is expected that the data written into the memory object be preserved on disk. A memory object may also apply to a physical device on an implementation. In this case, writes to the memory object are sent to the controller for the device and reads result in control registers being returned.

Mount Point*

The directory on which a “mounted file system” is mounted. This term, like *mount()* and *umount()*, was not included because it was implementation-defined.

Mounted File System*

See [File System](#) (on page 17).

834 **Name**

835 There are no explicit limits in IEEE Std 1003.1-200x on the sizes of names, words (see the
 836 definition of word in the Base Definitions volume of IEEE Std 1003.1-200x), lines, or other
 837 objects. However, other implicit limits do apply: shell script lines produced by many of the
 838 standard utilities cannot exceed {LINE_MAX} and the sum of exported variables comes under
 839 the {ARG_MAX} limit. Historical shells dynamically allocate memory for names and words and
 840 parse incoming lines a character at a time. Lines cannot have an arbitrary {LINE_MAX} limit
 841 because of historical practice, such as makefiles, where *make* removes the <newline>s associated
 842 with the commands for a target and presents the shell with one very long line. The text on
 843 INPUT FILES in the Shell and Utilities volume of IEEE Std 1003.1-200x, Section 1.11, Utility
 844 Description Defaults does allow a shell to run out of memory, but it cannot have arbitrary
 845 programming limits.

846 **Native Implementation***

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

850 **Nice Value**

851 This definition is not intended to suggest that all processes in a system have priorities that are
 852 comparable. Scheduling policy extensions, such as adding realtime priorities, make the notion of
 853 a single underlying priority for all scheduling policies problematic. Some implementations may
 854 implement the features related to *nice* to affect all processes on the system, others to affect just
 855 the general time-sharing activities implied by IEEE Std 1003.1-200x, and others may have no
 856 effect at all. Because of the use of “implementation-defined” in *nice* and *renice*, a wide range of
 857 implementation strategies is possible.

858 **Open File Description**

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

862 The following alternate names were discussed:

- 863 • For “open file description”:
 864 “open instance”, “file access description”, “open file information”, and “file access
 865 information”.
- 866 • For “file descriptor”:
 867 “file handle”, “file number” (cf., *fileno()*). Some historical implementations use the term
 868 “file table entry”.

869 **Orphaned Process Group**

870 Historical implementations have a concept of an orphaned process, which is a process whose
 871 parent process has exited. When job control is in use, it is necessary to prevent processes from
 872 being stopped in response to interactions with the terminal after they no longer are controlled by
 873 a job control-cognizant program. Because signals generated by the terminal are sent to a process
 874 group and not to individual processes, and because a signal may be provoked by a process that
 875 is not orphaned, but sent to another process that is orphaned, it is necessary to define an
 876 orphaned process group. The definition assumes that a process group will be manipulated as a
 877 group and that the job control-cognizant process controlling the group is outside of the group
 878 and is the parent of at least one process in the group (so that state changes may be reported via

879 *waitpid()*). Therefore, a group is considered to be controlled as long as at least one process in the
880 group has a parent that is outside of the process group, but within the session.

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

884 **Page**

885 The term "page" is defined to support the description of the behavior of memory mapping for
886 shared memory and memory mapped files, and the description of the behavior of process
887 memory locking. It is not intended to imply that shared memory/file mapping and memory
888 locking are applicable only to "paged" architectures. For the purposes of IEEE Std 1003.1-200x,
889 whatever the granularity on which an architecture supports mapping or locking, this is
890 considered to be a "page" . If an architecture cannot support the memory mapping or locking
891 functions specified by IEEE Std 1003.1-200x on any granularity, then these options will not be
892 implemented on the architecture.

893 **Passwd File***

894 Implementation-defined; see [User Database](#) (on page 32).

895 **Parent Directory**

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

902 **Pipe**

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

905 **Portable Filename Character Set**

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

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

918 Nothing in POSIX.1 is intended to preclude the use of extended characters where interchange is
919 not required or where mutual agreement is obtained. It has been suggested that in several places
920 "should" be used instead of "shall". Because (in the worst case) use of any character beyond the
921 portable filename character set would render the program or data not portable to all possible
922 systems, no extensions are permitted in this context.

923 **Process Lifetime**

924 IEEE Std 1003.1-2001/Cor 2-2004, item XBD/TC2/D6/5 is applied, adding *fork()*, *posix_spawn()*,
 925 *posix_spawnnp()*, and *vfork()* to the list of functions.

926 **Process Termination**

927 IEEE Std 1003.1-2001/Cor 2-2004, item XBD/TC2/D6/6 is applied, rewording the definition to
 928 address the “passive exit” on termination of the last thread or the *_Exit()* function.

929 **Regular File**

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

933 **Root Directory**

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

936 **Root File System***

937 Implementation-defined.

938 **Root of a File System***

939 Implementation-defined; see [Mount Point*](#) (on page 22).

940 **Signal**

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

943 **Superuser***

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

946 **Supplementary Group ID**

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

- 954 1. Reword the description of *setgid()* to permit it to change the supplementary group IDs to
 955 reflect the new effective group ID. A problem with this is that it adds more “may”s to the
 956 wording and does not address the portability problems of this optional behavior.
- 957 2. Mandate the inclusion of the effective group ID in the supplementary set (giving
 958 {NGROUPS_MAX} a minimum value of 1). This is the behavior of 4.4 BSD. In that
 959 system, the effective group ID is the first element of the array of supplementary group
 960 IDs (there is no separate copy stored, and changes to the effective group ID are made only
 961 in the supplementary group set). By convention, the initial value of the effective group ID

962 is duplicated elsewhere in the array so that the initial value is not lost when executing a
963 set-group-ID program.

- 964 3. Change the definition of supplementary group ID to exclude the effective group ID and
965 specify that the effective group ID does not change the set of supplementary group IDs.
966 This is the behavior of 4.2 BSD, 4.3 BSD, and System V Release 4.
- 967 4. Change the definition of supplementary group ID to exclude the effective group ID, and
968 require that *getgroups()* return the union of the effective group ID and the supplementary
969 group IDs.
- 970 5. Change the definition of {NGROUPS_MAX} to be one more than the number of
971 supplementary group IDs, so it continues to be the number of values returned by
972 *getgroups()* and existing applications continue to work. This alternative is effectively the
973 same as the second (and might actually have the same implementation).

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

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

987 Symbolic Link

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

998 Historical implementations were followed when determining which interfaces should apply to
999 symbolic links. Interfaces that historically followed symbolic links include *chmod()*, *link()*, and
1000 *utime()*. Interfaces that historically do not follow symbolic links include *chown()*, *lstat()*,
1001 *readlink()*, *rename()*, *remove()*, *rmdir()*, and *unlink()*. IEEE Std 1003.1-200x deviates from
1002 historical practice only in the case of *chown()*. Because there is no requirement for systems not
1003 supporting the XSI extension that there is an association of ownership with symbolic links, there
1004 was no interface in the base standard to change ownership. In addition, other implementations
1005 of symbolic links have modified *chown()* to follow symbolic links.

1006 In the case of symbolic links, IEEE Std 1003.1-200x states that a trailing slash is considered to be
1007 the final component of a pathname rather than the pathname component that preceded it. This
1008 is the behavior of historical implementations. For example, for */a/b* and */a/b/*, if */a/b* is a symbolic
1009 link to a directory, then */a/b* refers to the symbolic link, and */a/b/* is the same as */a/b/.*, which is

1010 the directory to which the symbolic link points.

1011 For multi-level security purposes, it is possible to have the link read mode govern permission
 1012 for the *readlink()* function. It is also possible that the read permissions of the directory containing
 1013 the link be used for this purpose. Implementations may choose to use either of these methods;
 1014 however, this is not current practice and neither method is specified.

1015 Several reasons were advanced for requiring that when a symbolic link is used as the source
 1016 argument to the *link()* function, the resulting link will apply to the file named by the contents of
 1017 the symbolic link rather than to the symbolic link itself. This is the case in historical
 1018 implementations. This action was preferred, as it supported the traditional idea of persistence
 1019 with respect to the target of a hard link. This decision is appropriate in light of a previous
 1020 decision not to require association of attributes with symbolic links, thereby allowing
 1021 implementations which do not use inodes. Opposition centered on the lack of symmetry on the
 1022 part of the *link()* and *unlink()* function pair with respect to symbolic links.

1023 Because a symbolic link and its referenced object coexist in the file system name space, confusion
 1024 can arise in distinguishing between the link itself and the referenced object. Historically, utilities
 1025 and system calls have adopted their own link following conventions in a somewhat *ad hoc*
 1026 fashion. Rules for a uniform approach are outlined here, although historical practice has been
 1027 adhered to as much as was possible. To promote consistent system use, user-written utilities are
 1028 encouraged to follow these same rules.

1029 Symbolic links are handled either by operating on the link itself, or by operating on the object
 1030 referenced by the link. In the latter case, an application or system call is said to “follow” the link.
 1031 Symbolic links may reference other symbolic links, in which case links are dereferenced until an
 1032 object that is not a symbolic link is found, a symbolic link that references a file that does not exist
 1033 is found, or a loop is detected. (Current implementations do not detect loops, but have a limit on
 1034 the number of symbolic links that they will dereference before declaring it an error.)

1035 There are four domains for which default symbolic link policy is established in a system. In
 1036 almost all cases, there are utility options that override this default behavior. The four domains
 1037 are as follows:

- 1038 1. Symbolic links specified to system calls that take filename arguments
- 1039 2. Symbolic links specified as command line filename arguments to utilities that are not
 1040 performing a traversal of a file hierarchy
- 1041 3. Symbolic links referencing files not of type directory, specified to utilities that are
 1042 performing a traversal of a file hierarchy
- 1043 4. Symbolic links referencing files of type directory, specified to utilities that are performing
 1044 a traversal of a file hierarchy

1045 *First Domain*

1046 The first domain is considered in earlier rationale.

1047 *Second Domain*

1048 The reason this category is restricted to utilities that are not traversing the file hierarchy is that
 1049 some standard utilities take an option that specifies a hierarchical traversal, but by default
 1050 operate on the arguments themselves. Generally, users specifying the option for a file hierarchy
 1051 traversal wish to operate on a single, physical hierarchy, and therefore symbolic links, which
 1052 may reference files outside of the hierarchy, are ignored. For example, *chown owner file* is a
 1053 different operation from the same command with the *-R* option specified. In this example, the
 1054 behavior of the command *chown owner file* is described here, while the behavior of the command
 1055 *chown -R owner file* is described in the third and fourth domains.

1056 The general rule is that the utilities in this category follow symbolic links named as arguments.

1057 Exceptions in the second domain are:

- 1058 • The *mv* and *rm* utilities do not follow symbolic links named as arguments, but respectively
- 1059 attempt to rename or delete them.
- 1060 • The *ls* utility is also an exception to this rule. For compatibility with historical systems,
- 1061 when the **-R** option is not specified, the *ls* utility follows symbolic links named as
- 1062 arguments if the **-L** option is specified or if the **-F**, **-d**, or **-l** options are not specified. (If
- 1063 the **-L** option is specified, *ls* always follows symbolic links; it is the only utility where the
- 1064 **-L** option affects its behavior even though a tree walk is not being performed.)

1065 All other standard utilities, when not traversing a file hierarchy, always follow symbolic links

1066 named as arguments.

1067 Historical practice is that the **-h** option is specified if standard utilities are to act upon symbolic

1068 links instead of upon their targets. Examples of commands that have historically had a **-h** option

1069 for this purpose are the *chgrp*, *chown*, *file*, and *test* utilities.

1070 *Third Domain*

1071 The third domain is symbolic links, referencing files not of type directory, specified to utilities

1072 that are performing a traversal of a file hierarchy. (This includes symbolic links specified as

1073 command line filename arguments or encountered during the traversal.)

1074 The intention of the Shell and Utilities volume of IEEE Std 1003.1-200x is that the operation that

1075 the utility is performing is applied to the symbolic link itself, if that operation is applicable to

1076 symbolic links. The reason that the operation is not required is that symbolic links in some

1077 implementations do not have such attributes as a file owner, and therefore the *chown* operation

1078 would be meaningless. If symbolic links on the system have an owner, it is the intention that the

1079 utility *chown* cause the owner of the symbolic link to change. If symbolic links do not have an

1080 owner, the symbolic link should be ignored. Specifically, by default, no change should be made

1081 to the file referenced by the symbolic link.

1082 *Fourth Domain*

1083 The fourth domain is symbolic links referencing files of type directory, specified to utilities that

1084 are performing a traversal of a file hierarchy. (This includes symbolic links specified as

1085 command line filename arguments or encountered during the traversal.)

1086 Most standard utilities do not, by default, indirect into the file hierarchy referenced by the

1087 symbolic link. (The Shell and Utilities volume of IEEE Std 1003.1-200x uses the informal term

1088 “physical walk” to describe this case. The case where the utility does indirect through the

1089 symbolic link is termed a “logical walk”.)

1090 There are three reasons for the default to be a physical walk:

- 1091 1. With very few exceptions, a physical walk has been the historical default on UNIX
- 1092 systems supporting symbolic links. Because some utilities (that is, *rm*) must default to a
- 1093 physical walk, regardless, changing historical practice in this regard would be confusing
- 1094 to users and needlessly incompatible.
- 1095 2. For systems where symbolic links have the historical file attributes (that is, *owner*, *group*,
- 1096 *mode*), defaulting to a logical traversal would require the addition of a new option to the
- 1097 commands to modify the attributes of the link itself. This is painful and more complex
- 1098 than the alternatives.

- 1099 3. There is a security issue with defaulting to a logical walk. Historically, the command
 1100 *chown -R user file* has been safe for the superuser because *setuid* and *setgid* bits were lost
 1101 when the ownership of the file was changed. If the walk were logical, changing
 1102 ownership would no longer be safe because a user might have inserted a symbolic link
 1103 pointing to any file in the tree. Again, this would necessitate the addition of an option to
 1104 the commands doing hierarchy traversal to not indirect through the symbolic links, and
 1105 historical scripts doing recursive walks would instantly become security problems. While
 1106 this is mostly an issue for system administrators, it is preferable to not have different
 1107 defaults for different classes of users.

1108 However, the standard developers agreed to leave it unspecified to achieve consensus.

1109 As consistently as possible, users may cause standard utilities performing a file hierarchy
 1110 traversal to follow any symbolic links named on the command line, regardless of the type of file
 1111 they reference, by specifying the **-H** (for half logical) option. This option is intended to make the
 1112 command line name space look like the logical name space.

1113 As consistently as possible, users may cause standard utilities performing a file hierarchy
 1114 traversal to follow any symbolic links named on the command line as well as any symbolic links
 1115 encountered during the traversal, regardless of the type of file they reference, by specifying the
 1116 **-L** (for logical) option. This option is intended to make the entire name space look like the
 1117 logical name space.

1118 For consistency, implementors are encouraged to use the **-P** (for “physical”) flag to specify the
 1119 physical walk in utilities that do logical walks by default for whatever reason. The only standard
 1120 utilities that require the **-P** option are *cd* and *pwd*; see the note below.

1121 When one or more of the **-H**, **-L**, and **-P** flags can be specified, the last one specified determines
 1122 the behavior of the utility. This permits users to alias commands so that the default behavior is a
 1123 logical walk and then override that behavior on the command line.

1124 *Exceptions in the Third and Fourth Domains*

1125 The *ls* and *rm* utilities are exceptions to these rules. The *rm* utility never follows symbolic links
 1126 and does not support the **-H**, **-L**, or **-P** options. Some historical versions of *ls* always followed
 1127 symbolic links given on the command line whether the **-L** option was specified or not.
 1128 Historical versions of *ls* did not support the **-H** option. In IEEE Std 1003.1-200x, unless one of
 1129 the **-H** or **-L** options is specified, the *ls* utility only follows symbolic links to directories that are
 1130 given as operands. The *ls* utility does not support the **-P** option.

1131 The Shell and Utilities volume of IEEE Std 1003.1-200x requires that the standard utilities *ls*, *find*,
 1132 and *pax* detect infinite loops when doing logical walks; that is, a directory, or more commonly a
 1133 symbolic link, that refers to an ancestor in the current file hierarchy. If the file system itself is
 1134 corrupted, causing the infinite loop, it may be impossible to recover. Because *find* and *ls* are often
 1135 used in system administration and security applications, they should attempt to recover and
 1136 continue as best as they can. The *pax* utility should terminate because the archive it was creating
 1137 is by definition corrupted. Other, less vital, utilities should probably simply terminate as well.
 1138 Implementations are strongly encouraged to detect infinite loops in all utilities.

1139 Historical practice is shown in [Table A-1](#) (on page 30). The heading **SVID3** stands for the Third
 1140 Edition of the System V Interface Definition.

1141 Historically, several shells have had built-in versions of the *pwd* utility. In some of these shells,
 1142 *pwd* reported the physical path, and in others, the logical path. Implementations of the shell
 1143 corresponding to IEEE Std 1003.1-200x must report the logical path by default. Earlier versions
 1144 of IEEE Std 1003.1-200x did not require the *pwd* utility to be a built-in utility. Now that *pwd*
 1145 is required to set an environment variable in the current shell execution environment, it must be a
 1146 built-in utility.

1147 The *cd* command is required, by default, to treat the filename dot-dot logically. Implementors are
 1148 required to support the **-P** flag in *cd* so that users can have their current environment handled
 1149 physically. In 4.3 BSD, *chgrp* during tree traversal changed the group of the symbolic link, not
 1150 the target. Symbolic links in 4.4 BSD do not have *owner*, *group*, *mode*, or other standard UNIX
 1151 system file attributes.

1152 **Table A-1** Historical Practice for Symbolic Links

Utility	SVID3	4.3 BSD	4.4 BSD	POSIX	Comments
1153 <i>cd</i>				-L	Treat ". ." logically.
1154 <i>cd</i>				-P	Treat ". ." physically.
1155 <i>chgrp</i>			-H	-H	Follow command line symlinks.
1156 <i>chgrp</i>			-h	-L	Follow symlinks.
1157 <i>chgrp</i>	-h			-h	Affect the symlink.
1158 <i>chmod</i>					Affect the symlink.
1159 <i>chmod</i>			-H		Follow command line symlinks.
1160 <i>chmod</i>			-h		Follow symlinks.
1161 <i>chown</i>			-H	-H	Follow command line symlinks.
1162 <i>chown</i>			-h	-L	Follow symlinks.
1163 <i>chown</i>	-h			-h	Affect the symlink.
1164 <i>cp</i>			-H	-H	Follow command line symlinks.
1165 <i>cp</i>			-h	-L	Follow symlinks.
1166 <i>cpio</i>	-L		-L		Follow symlinks.
1167 <i>du</i>			-H	-H	Follow command line symlinks.
1168 <i>du</i>			-h	-L	Follow symlinks.
1169 <i>file</i>	-h			-h	Affect the symlink.
1170 <i>find</i>			-H	-H	Follow command line symlinks.
1171 <i>find</i>			-h	-L	Follow symlinks.
1172 <i>find</i>	-follow		-follow		Follow symlinks.
1173 <i>ln</i>	-s	-s	-s	-s	Create a symbolic link.
1174 <i>ls</i>	-L	-L	-L	-L	Follow symlinks.
1175 <i>ls</i>				-H	Follow command line symlinks.
1176 <i>mv</i>					Operates on the symlink.
1177 <i>pax</i>			-H	-H	Follow command line symlinks.
1178 <i>pax</i>			-h	-L	Follow symlinks.
1179 <i>pwd</i>				-L	Printed path may contain symlinks.
1180 <i>pwd</i>				-P	Printed path will not contain symlinks.
1181 <i>rm</i>					Operates on the symlink.
1182 <i>tar</i>			-H		Follow command line symlinks.
1183 <i>tar</i>		-h	-h		Follow symlinks.
1184 <i>test</i>	-h		-h	-h	Affect the symlink.

1186 Synchronously-Generated Signal

1187 Those signals that may be generated synchronously include SIGABRT, SIGBUS, SIGILL, SIGFPE,
 1188 SIGPIPE, and SIGSEGV.

1189 Any signal sent via the *raise()* function or a *kill()* function targeting the current process is also
 1190 considered synchronous.

1191 **System Call***

1192 The distinction between a “system call” and a “library routine” is an implementation detail that
 1193 may differ between implementations and has thus been excluded from POSIX.1.

1194 See “Interface, Not Implementation” in [Introduction](#) (on page 0).

1195 **System Console**

1196 IEEE Std 1003.1-2001/Cor 2-2004, item XBD/TC2/D6/7 is applied, changing from “An
 1197 implementation-defined device” to “A device”.

1198 **System Databases**

1199 IEEE Std 1003.1-2001/Cor 2-2004, item XBD/TC2/D6/9 is applied, rewording the definition to
 1200 reference the existing definitions for “group database” and “user database”.

1201 **System Process**

1202 IEEE Std 1003.1-2001/Cor 2-2004, item XBD/TC2/D6/8 is applied, rewording the definition to
 1203 remove the requirement for an implementation to define the object.

1204 **System Reboot**

1205 A “system reboot” is an event initiated by an unspecified circumstance that causes all processes
 1206 (other than special system processes) to be terminated in an implementation-defined manner,
 1207 after which any changes to the state and contents of files created or written to by a Conforming
 1208 POSIX.1 Application prior to the event are implementation-defined.

1209 IEEE Std 1003.1-2001/Cor 2-2004, item XBD/TC2/D6/10 is applied, changing “An
 1210 implementation-defined sequence of events” to “An unspecified sequence of events”.

1211 **Synchronized I/O Data (and File) Integrity Completion**

1212 These terms specify that for synchronized read operations, pending writes must be successfully
 1213 completed before the read operation can complete. This is motivated by two circumstances.
 1214 Firstly, when synchronizing processes can access the same file, but not share common buffers
 1215 (such as for a remote file system), this requirement permits the reading process to guarantee that
 1216 it can read data written remotely. Secondly, having data written synchronously is insufficient to
 1217 guarantee the order with respect to a subsequent write by a reading process, and thus this extra
 1218 read semantic is necessary.

1219 **Text File**

1220 The term “text file” does not prevent the inclusion of control or other non-printable characters
 1221 (other than NUL). Therefore, standard utilities that list text files as inputs or outputs are either
 1222 able to process the special characters or they explicitly describe their limitations within their
 1223 individual descriptions. The definition of “text file” has caused controversy. The only difference
 1224 between text and binary files is that text files have lines of less than {LINE_MAX} bytes, with no
 1225 NUL characters, each terminated by a <newline>. The definition allows a file with a single
 1226 <newline>, but not a totally empty file, to be called a text file. If a file ends with an incomplete
 1227 line it is not strictly a text file by this definition. The <newline> referred to in
 1228 IEEE Std 1003.1-200x is not some generic line separator, but a single character; files created on
 1229 systems where they use multiple characters for ends of lines are not portable to all conforming
 1230 systems without some translation process unspecified by IEEE Std 1003.1-200x.

1231 **Thread**

1232 IEEE Std 1003.1-200x defines a thread to be a flow of control within a process. Each thread has a
 1233 minimal amount of private state; most of the state associated with a process is shared among all
 1234 of the threads in the process. While most multi-thread extensions to POSIX have taken this
 1235 approach, others have made different decisions.

1236 **Note:** The choice to put threads within a process does not constrain implementations to implement
 1237 threads in that manner. However, all functions have to behave as though threads share the
 1238 indicated state information with the process from which they were created.

1239 Threads need to share resources in order to cooperate. Memory has to be widely shared between
 1240 threads in order for the threads to cooperate at a fine level of granularity. Threads keep data
 1241 structures and the locks protecting those data structures in shared memory. For a data structure
 1242 to be usefully shared between threads, such structures should not refer to any data that can only
 1243 be interpreted meaningfully by a single thread. Thus, any system resources that might be
 1244 referred to in data structures need to be shared between all threads. File descriptors, pathnames,
 1245 and pointers to stack variables are all things that programmers want to share between their
 1246 threads. Thus, the file descriptor table, the root directory, the current working directory, and the
 1247 address space have to be shared.

1248 Library implementations are possible as long as the effective behavior is as if system services
 1249 invoked by one thread do not suspend other threads. This may be difficult for some library
 1250 implementations on systems that do not provide asynchronous facilities.

1251 See [Section B.2.9](#) for additional rationale.

1252 **Thread ID**

1253 See [Section B.2.9.2](#) for additional rationale.

1254 **Thread-Safe Function**

1255 All functions required by IEEE Std 1003.1-200x need to be thread-safe; see [Section A.4.17](#) and
 1256 [Section B.2.9.1](#) for additional rationale.

1257 **User Database**

1258 There are no references in IEEE Std 1003.1-200x to a “passwd file” or a “group file”, and there is
 1259 no requirement that the *group* or *passwd* databases be kept in files containing editable text. Many
 1260 large timesharing systems use *passwd* databases that are hashed for speed. Certain security
 1261 classifications prohibit certain information in the *passwd* database from being publicly readable.

1262 The term “encoded” is used instead of “encrypted” in order to avoid the implementation
 1263 connotations (such as reversibility or use of a particular algorithm) of the latter term.

1264 The *getgrent()*, *setgrent()*, *endgrent()*, *getpwent()*, *setpwent()*, and *endpwent()* functions are not
 1265 included as part of the base standard because they provide a linear database search capability
 1266 that is not generally useful (the *getpwuid()*, *getpwnam()*, *getgrgid()*, and *getgrnam()* functions are
 1267 provided for keyed lookup) and because in certain distributed systems, especially those with
 1268 different authentication domains, it may not be possible or desirable to provide an application
 1269 with the ability to browse the system databases indiscriminately. They are provided on XSI-
 1270 conformant systems due to their historical usage by many existing applications.

1271 A change from historical implementations is that the structures used by these functions have
 1272 fields of the types **gid_t** and **uid_t**, which are required to be defined in the **<sys/types.h>** header.
 1273 IEEE Std 1003.1-200x requires implementations to ensure that these types are defined by
 1274 inclusion of **<grp.h>** and **<pwd.h>**, respectively, without imposing any name space pollution or
 1275 errors from redefinition of types.

1276 IEEE Std 1003.1-200x is silent about the content of the strings containing user or group names.
 1277 These could be digit strings. IEEE Std 1003.1-200x is also silent as to whether such digit strings
 1278 bear any relationship to the corresponding (numeric) user or group ID.

1279 *Database Access*

1280 The thread-safe versions of the user and group database access functions return values in user-
 1281 supplied buffers instead of possibly using static data areas that may be overwritten by each call.

1282 IEEE Std 1003.1-2001/Cor 2-2004, item XBD/TC2/D6/11 is applied, removing the words “of
 1283 implementation-defined format”.

1284 **Virtual Processor***

1285 The term “virtual processor” was chosen as a neutral term describing all kernel-level
 1286 schedulable entities, such as processes, Mach tasks, or lightweight processes. Implementing
 1287 threads using multiple processes as virtual processors, or implementing multiplexed threads
 1288 above a virtual processor layer, should be possible, provided some mechanism has also been
 1289 implemented for sharing state between processes or virtual processors. Many systems may also
 1290 wish to provide implementations of threads on systems providing “shared processes” or
 1291 “variable-weight processes”. It was felt that exposing such implementation details would
 1292 severely limit the type of systems upon which the threads interface could be supported and
 1293 prevent certain types of valid implementations. It was also determined that a virtual processor
 1294 interface was out of the scope of the Rationale (Informative) volume of IEEE Std 1003.1-200x.

1295 **XSI**

1296 This is included to allow IEEE Std 1003.1-200x to be adopted as an IEEE standard and an Open
 1297 Group Technical Standard, serving both the POSIX and the Single UNIX Specification in a core
 1298 set of volumes.

1299 The term “XSI” has been used for 10 years in connection with the XPG series and the first and
 1300 second versions of the base volumes of the Single UNIX Specification. The XSI margin code was
 1301 introduced to denote the extended or more restrictive semantics beyond POSIX that are
 1302 applicable to UNIX systems.

1303 **A.4 General Concepts**

1304 The general concepts are similar in nature to the definitions section, with the exception that a
 1305 term defined in general concepts can contain normative requirements.

1306 **A.4.1 Concurrent Execution**

1307 There is no additional rationale provided for this section.

1308 **A.4.2 Directory Protection**

1309 There is no additional rationale provided for this section.

1310 A.4.3 Extended Security Controls

1311 Allowing an implementation to define extended security controls enables the use of
 1312 IEEE Std 1003.1-200x in environments that require different or more rigorous security than that
 1313 provided in POSIX.1. Extensions are allowed in two areas: privilege and file access permissions.
 1314 The semantics of these areas have been defined to permit extensions with reasonable, but not
 1315 exact, compatibility with all existing practices. For example, the elimination of the superuser
 1316 definition precludes identifying a process as privileged or not by virtue of its effective user ID.

1317 A.4.4 File Access Permissions

1318 A process should not try to anticipate the result of an attempt to access data by *a priori* use of
 1319 these rules. Rather, it should make the attempt to access data and examine the return value (and
 1320 possibly *errno* as well), or use *access()*. An implementation may include other security
 1321 mechanisms in addition to those specified in POSIX.1, and an access attempt may fail because of
 1322 those additional mechanisms, even though it would succeed according to the rules given in this
 1323 section. (For example, the user's security level might be lower than that of the object of the
 1324 access attempt.) The supplementary group IDs provide another reason for a process to not
 1325 attempt to anticipate the result of an access attempt.

1326 Since the current standard does not specify a method for opening a directory for searching, it is
 1327 unspecified whether search permission on the *fd* argument to *openat()* and related functions is
 1328 based on whether the directory was opened with search mode or on the current permissions
 1329 allowed by the directory at the time a search is performed. When there is existing practice that
 1330 supports opening directories for searching, it is expected that these functions will be modified to
 1331 specify that the search permissions will be granted based on the file access modes of the
 1332 directory's file descriptor identified by *fd*, and not on the mode of the directory at the time the
 1333 directory is searched.

1334 A.4.5 File Hierarchy

1335 Though the file hierarchy is commonly regarded to be a tree, POSIX.1 does not define it as such
 1336 for three reasons:

- 1337 1. Links may join branches.
- 1338 2. In some network implementations, there may be no single absolute root directory; see
 1339 *pathname resolution*.
- 1340 3. With symbolic links, the file system need not be a tree or even a directed acyclic graph.

1341 A.4.6 Filenames

1342 Historically, certain filenames have been reserved. This list includes **core**, **/etc/passwd**, and so
 1343 on. Conforming applications should avoid these.

1344 Most historical implementations prohibit case folding in filenames; that is, treating uppercase
 1345 and lowercase alphabetic characters as identical. However, some consider case folding desirable:

- 1346 • For user convenience
- 1347 • For ease-of-implementation of the POSIX.1 interface as a hosted system on some popular
 1348 operating systems

1349 Variants, such as maintaining case distinctions in filenames, but ignoring them in comparisons,
 1350 have been suggested. Methods of allowing escaped characters of the case opposite the default
 1351 have been proposed.

1352 Many reasons have been expressed for not allowing case folding, including:

- 1353
- 1354
- 1355
- 1356
- 1357
- 1358
- 1359
- 1360
- 1361
- 1362
- 1363
- 1364
- 1365
- 1366
- 1367
- 1368
- 1369
- 1370
- 1371
- 1372
- 1373
- 1374
- 1375
- 1376
- 1377
- 1378
- 1379
- 1380
- 1381
- 1382
- 1383
- 1384
- 1385
- 1386
- 1387
- 1388
- 1389
- No solid evidence has been produced as to whether case-sensitivity or case-insensitivity is more convenient for users.
 - Making case-insensitivity a POSIX.1 implementation option would be worse than either having it or not having it, because:
 - More confusion would be caused among users.
 - Application developers would have to account for both cases in their code.
 - POSIX.1 implementors would still have other problems with native file systems, such as short or otherwise constrained filenames or pathnames, and the lack of hierarchical directory structure.
 - Case folding is not easily defined in many European languages, both because many of them use characters outside the US ASCII alphabetic set, and because:
 - In Spanish, the digraph "ll" is considered to be a single letter, the capitalized form of which may be either "Ll" or "LL", depending on context.
 - In French, the capitalized form of a letter with an accent may or may not retain the accent, depending on the country in which it is written.
 - In German, the sharp ess may be represented as a single character resembling a Greek beta (β) in lowercase, but as the digraph "SS" in uppercase.
 - In Greek, there are several lowercase forms of some letters; the one to use depends on its position in the word. Arabic has similar rules.
 - Many East Asian languages, including Japanese, Chinese, and Korean, do not distinguish case and are sometimes encoded in character sets that use more than one byte per character.
 - Multiple character codes may be used on the same machine simultaneously. There are several ISO character sets for European alphabets. In Japan, several Japanese character codes are commonly used together, sometimes even in filenames; this is evidently also the case in China. To handle case insensitivity, the kernel would have to at least be able to distinguish for which character sets the concept made sense.
 - The file system implementation historically deals only with bytes, not with characters, except for slash and the null byte.
 - The purpose of POSIX.1 is to standardize the common, existing definition, not to change it. Mandating case-insensitivity would make all historical implementations non-standard.
 - Not only the interface, but also application programs would need to change, counter to the purpose of having minimal changes to existing application code.
 - At least one of the original developers of the UNIX system has expressed objection in the strongest terms to either requiring case-insensitivity or making it an option, mostly on the basis that POSIX.1 should not hinder portability of application programs across related implementations in order to allow compatibility with unrelated operating systems.

1390 Two proposals were entertained regarding case folding in filenames:

- 1391 1. Remove all wording that previously permitted case folding.

1392 Rationale Case folding is inconsistent with portable filename character set definition

1393 and filename definition (all characters except slash and null). No known

1394 implementations allowing all characters except slash and null also do case

1395 folding.

1396 2. Change “though this practice is not recommended:” to “although this practice is strongly
1397 discouraged.”

1398 Rationale If case folding must be included in POSIX.1, the wording should be stronger
1399 to discourage the practice.

1400 The consensus selected the first proposal. Otherwise, a conforming application would have to
1401 assume that case folding would occur when it was not wanted, but that it would not occur when
1402 it was wanted.

1403 **A.4.7 Filename Portability**

1404 Filenames should be constructed from the portable filename character set because the use of
1405 other characters can be confusing or ambiguous in certain contexts. (For example, the use of a
1406 colon (‘ : ’) in a pathname could cause ambiguity if that pathname were included in a *PATH*
1407 definition.)

1408 The constraint on use of the hyphen character as the first character of a portable filename is a
1409 constraint on application behavior and not on implementations, since applications might not
1410 work as expected when such a filename is passed as a command line argument.

1411 **A.4.8 File Times Update**

1412 This section reflects the actions of historical implementations. The times are not updated
1413 immediately, but are only marked for update by the functions. An implementation may update
1414 these times immediately.

1415 The accuracy of the time update values is intentionally left unspecified so that systems can
1416 control the bandwidth of a possible covert channel.

1417 The wording was carefully chosen to make it clear that there is no requirement that the
1418 conformance document contain information that might incidentally affect file update times. Any
1419 function that performs pathname resolution might update several *st_atime* fields. Functions such
1420 as *getpwnam()* and *getgrnam()* might update the *st_atime* field of some specific file or files. It is
1421 intended that these are not required to be documented in the conformance document, but they
1422 should appear in the system documentation.

1423 **A.4.9 Host and Network Byte Order**

1424 There is no additional rationale provided for this section.

1425 **A.4.10 Measurement of Execution Time**

1426 The methods used to measure the execution time of processes and threads, and the precision of
1427 these measurements, may vary considerably depending on the software architecture of the
1428 implementation, and on the underlying hardware. Implementations can also make tradeoffs
1429 between the scheduling overhead and the precision of the execution time measurements.
1430 IEEE Std 1003.1-200x does not impose any requirement on the accuracy of the execution time; it
1431 instead specifies that the measurement mechanism and its precision are implementation-
1432 defined.

1433 A.4.11 Memory Synchronization

1434 In older multi-processors, access to memory by the processors was strictly multiplexed. This
 1435 meant that a processor executing program code interrogates or modifies memory in the order
 1436 specified by the code and that all the memory operation of all the processors in the system
 1437 appear to happen in some global order, though the operation histories of different processors are
 1438 interleaved arbitrarily. The memory operations of such machines are said to be sequentially
 1439 consistent. In this environment, threads can synchronize using ordinary memory operations. For
 1440 example, a producer thread and a consumer thread can synchronize access to a circular data
 1441 buffer as follows:

```

1442 int rdptr = 0;
1443 int wrptr = 0;
1444 data_t buf[BUFSIZE];

1445 Thread 1:
1446     while (work_to_do) {
1447         int next;
1448
1449         buf[wrptr] = produce();
1450         next = (wrptr + 1) % BUFSIZE;
1451         while (rdptr == next)
1452             ;
1453         wrptr = next;
1454     }
1455
1456 Thread 2:
1457     while (work_to_do) {
1458         while (rdptr == wrptr)
1459             ;
1460         consume(buf[rdptr]);
1461         rdptr = (rdptr + 1) % BUFSIZE;
1462     }
  
```

1461 In modern multi-processors, these conditions are relaxed to achieve greater performance. If one
 1462 processor stores values in location A and then location B, then other processors loading data
 1463 from location B and then location A may see the new value of B but the old value of A. The
 1464 memory operations of such machines are said to be weakly ordered. On these machines, the
 1465 circular buffer technique shown in the example will fail because the consumer may see the new
 1466 value of *wrptr* but the old value of the data in the buffer. In such machines, synchronization can
 1467 only be achieved through the use of special instructions that enforce an order on memory
 1468 operations. Most high-level language compilers only generate ordinary memory operations to
 1469 take advantage of the increased performance. They usually cannot determine when memory
 1470 operation order is important and generate the special ordering instructions. Instead, they rely on
 1471 the programmer to use synchronization primitives correctly to ensure that modifications to a
 1472 location in memory are ordered with respect to modifications and/or access to the same location
 1473 in other threads. Access to read-only data need not be synchronized. The resulting program is
 1474 said to be data race-free.

1475 Synchronization is still important even when accessing a single primitive variable (for example,
 1476 an integer). On machines where the integer may not be aligned to the bus data width or be
 1477 larger than the data width, a single memory load may require multiple memory cycles. This
 1478 means that it may be possible for some parts of the integer to have an old value while other
 1479 parts have a newer value. On some processor architectures this cannot happen, but portable
 1480 programs cannot rely on this.

1481 In summary, a portable multi-threaded program, or a multi-process program that shares

1482 writable memory between processes, has to use the synchronization primitives to synchronize
 1483 data access. It cannot rely on modifications to memory being observed by other threads in the
 1484 order written in the application or even on modification of a single variable being seen
 1485 atomically.

1486 Conforming applications may only use the functions listed to synchronize threads of control
 1487 with respect to memory access. There are many other candidates for functions that might also be
 1488 used. Examples are: signal sending and reception, or pipe writing and reading. In general, any
 1489 function that allows one thread of control to wait for an action caused by another thread of
 1490 control is a candidate. IEEE Std 1003.1-200x does not require these additional functions to
 1491 synchronize memory access since this would imply the following:

- 1492 • All these functions would have to be recognized by advanced compilation systems so that
- 1493 memory operations and calls to these functions are not reordered by optimization.
- 1494 • All these functions would potentially have to have memory synchronization instructions
- 1495 added, depending on the particular machine.
- 1496 • The additional functions complicate the model of how memory is synchronized and make
- 1497 automatic data race detection techniques impractical.

1498 Formal definitions of the memory model were rejected as unreadable by the vast majority of
 1499 programmers. In addition, most of the formal work in the literature has concentrated on the
 1500 memory as provided by the hardware as opposed to the application programmer through the
 1501 compiler and runtime system. It was believed that a simple statement intuitive to most
 1502 programmers would be most effective. IEEE Std 1003.1-200x defines functions that can be used
 1503 to synchronize access to memory, but it leaves open exactly how one relates those functions to
 1504 the semantics of each function as specified elsewhere in IEEE Std 1003.1-200x.
 1505 IEEE Std 1003.1-200x also does not make a formal specification of the partial ordering in time
 1506 that the functions can impose, as that is implied in the description of the semantics of each
 1507 function. It simply states that the programmer has to ensure that modifications do not occur
 1508 “simultaneously” with other access to a memory location.

1509 IEEE Std 1003.1-2001/Cor 1-2002, item XBD/TC1/D6/4 is applied, adding a new paragraph
 1510 beneath the table of functions: “The *pthread_once()* function shall synchronize memory for the
 1511 first call in each thread for a given **pthread_once_t** object.”.

1512 A.4.12 Pathname Resolution

1513 It is necessary to differentiate between the definition of pathname and the concept of pathname
 1514 resolution with respect to the handling of trailing slashes. By specifying the behavior here, it is
 1515 not possible to provide an implementation that is conforming but extends all interfaces that
 1516 handle pathnames to also handle strings that are not legal pathnames (because they have
 1517 trailing slashes).

1518 Pathnames that end with one or more trailing slash characters must refer to directory paths.
 1519 Previous versions of IEEE Std 1003.1-200x were not specific about the distinction between
 1520 trailing slashes on files and directories, and both were permitted.

1521 Two types of implementation have been prevalent; those that ignored trailing slash characters
 1522 on all pathnames regardless, and those that permitted them only on existing directories.

1523 IEEE Std 1003.1-200x requires that a pathname with a trailing slash character be treated as if it
 1524 had a trailing “/.” everywhere.

1525 Note that this change does not break any conforming applications; since there were two different
 1526 types of implementation, no application could have portably depended on either behavior. This
 1527 change does however require some implementations to be altered to remain compliant.
 1528 Substantial discussion over a three-year period has shown that the benefits to application

1529 developers outweighs the disadvantages for some vendors.

1530 On a historical note, some early applications automatically appended a `'/'` to every path.
 1531 Rather than fix the applications, the system implementation was modified to accept this
 1532 behavior by ignoring any trailing slash.

1533 Each directory has exactly one parent directory which is represented by the name **dot-dot** in the
 1534 first directory. No other directory, regardless of linkages established by symbolic links, is
 1535 considered the parent directory by IEEE Std 1003.1-200x.

1536 There are two general categories of interfaces involving pathname resolution: those that follow
 1537 the symbolic link, and those that do not. There are several exceptions to this rule; for example,
 1538 `open(path,O_CREAT|O_EXCL)` will fail when `path` names a symbolic link. However, in all other
 1539 situations, the `open()` function will follow the link.

1540 What the filename **dot-dot** refers to relative to the root directory is implementation-defined. In
 1541 Version 7 it refers to the root directory itself; this is the behavior mentioned in
 1542 IEEE Std 1003.1-200x. In some networked systems the construction `././hostname/` is used to refer
 1543 to the root directory of another host, and POSIX.1 permits this behavior.

1544 Other networked systems use the construct `//hostname` for the same purpose; that is, a double
 1545 initial slash is used. There is a potential problem with existing applications that create full
 1546 pathnames by taking a trunk and a relative pathname and making them into a single string
 1547 separated by `'/'`, because they can accidentally create networked pathnames when the trunk is
 1548 `'/'`. This practice is not prohibited because such applications can be made to conform by
 1549 simply changing to use `"/"` as a separator instead of `'/'`:

- 1550 • If the trunk is `'/'`, the full pathname will begin with `"/"` (the initial `'/'` and the
 1551 separator `"/"`). This is the same as `'/'`, which is what is desired. (This is the general
 1552 case of making a relative pathname into an absolute one by prefixing with `"/"` instead
 1553 of `'/'`.)
- 1554 • If the trunk is `"/A"`, the result is `"/A/..."`; since non-leading sequences of two or more
 1555 slashes are treated as a single slash, this is equivalent to the desired `"/A/..."`.
- 1556 • If the trunk is `"/A"`, the implementation-defined semantics will apply. (The multiple
 1557 slash rule would apply.)

1558 Application developers should avoid generating pathnames that start with `"/"`.
 1559 Implementations are strongly encouraged to avoid using this special interpretation since a
 1560 number of applications currently do not follow this practice and may inadvertently generate
 1561 `"/..."`.

1562 The term “root directory” is only defined in POSIX.1 relative to the process. In some
 1563 implementations, there may be no absolute root directory. The initialization of the root directory
 1564 of a process is implementation-defined.

1565 A.4.13 Process ID Reuse

1566 There is no additional rationale provided for this section.

1567 **A.4.14 Scheduling Policy**
 1568 There is no additional rationale provided for this section.

1569 **A.4.15 Seconds Since the Epoch**

1570 Coordinated Universal Time (UTC) includes leap seconds. However, in POSIX time (seconds
 1571 since the Epoch), leap seconds are ignored (not applied) to provide an easy and compatible
 1572 method of computing time differences. Broken-down POSIX time is therefore not necessarily
 1573 UTC, despite its appearance.

1574 As of September 2000, 24 leap seconds had been added to UTC since the Epoch, 1 January, 1970.
 1575 Historically, one leap second is added every 15 months on average, so this offset can be expected
 1576 to grow steadily with time.

1577 Most systems' notion of "time" is that of a continuously increasing value, so this value should
 1578 increase even during leap seconds. However, not only do most systems not keep track of leap
 1579 seconds, but most systems are probably not synchronized to any standard time reference.
 1580 Therefore, it is inappropriate to require that a time represented as seconds since the Epoch
 1581 precisely represent the number of seconds between the referenced time and the Epoch.

1582 It is sufficient to require that applications be allowed to treat this time as if it represented the
 1583 number of seconds between the referenced time and the Epoch. It is the responsibility of the
 1584 vendor of the system, and the administrator of the system, to ensure that this value represents
 1585 the number of seconds between the referenced time and the Epoch as closely as necessary for the
 1586 application being run on that system.

1587 It is important that the interpretation of time names and seconds since the Epoch values be
 1588 consistent across conforming systems; that is, it is important that all conforming systems
 1589 interpret "536 457 599 seconds since the Epoch" as 59 seconds, 59 minutes, 23 hours 31 December
 1590 1986, regardless of the accuracy of the system's idea of the current time. The expression is given
 1591 to ensure a consistent interpretation, not to attempt to specify the calendar. The relationship
 1592 between *tm_yday* and the day of week, day of month, and month is in accordance with the
 1593 Gregorian calendar, and so is not specified in POSIX.1.

1594 Consistent interpretation of seconds since the Epoch can be critical to certain types of distributed
 1595 applications that rely on such timestamps to synchronize events. The accrual of leap seconds in a
 1596 time standard is not predictable. The number of leap seconds since the Epoch will likely
 1597 increase. POSIX.1 is more concerned about the synchronization of time between applications of
 1598 astronomically short duration.

1599 Note that *tm_yday* is zero-based, not one-based, so the day number in the example above is 364.
 1600 Note also that the division is an integer division (discarding remainder) as in the C language.

1601 Note also that the meaning of *gmtime()*, *localtime()*, and *mktime()* is specified in terms of this
 1602 expression. However, the ISO C standard computes *tm_yday* from *tm_mday*, *tm_mon*, and *tm_year*
 1603 in *mktime()*. Because it is stated as a (bidirectional) relationship, not a function, and because the
 1604 conversion between month-day-year and day-of-year dates is presumed well known and is also
 1605 a relationship, this is not a problem.

1606 Implementations that implement **time_t** as a signed 32-bit integer will overflow in 2038. The
 1607 data size for **time_t** is as per the ISO C standard definition, which is implementation-defined.

1608 See also [Epoch](#) (on page 17).

1609 The topic of whether seconds since the Epoch should account for leap seconds has been debated
 1610 on a number of occasions, and each time consensus was reached (with acknowledged dissent
 1611 each time) that the majority of users are best served by treating all days identically. (That is, the
 1612 majority of applications were judged to assume a single length—as measured in seconds since

1613 the Epoch—for all days. Thus, leap seconds are not applied to seconds since the Epoch.) Those
 1614 applications which do care about leap seconds can determine how to handle them in whatever
 1615 way those applications feel is best. This was particularly emphasized because there was
 1616 disagreement about what the best way of handling leap seconds might be. It is a practical
 1617 impossibility to mandate that a conforming implementation must have a fixed relationship to
 1618 any particular official clock (consider isolated systems, or systems performing “reruns” by
 1619 setting the clock to some arbitrary time).

1620 Note that as a practical consequence of this, the length of a second as measured by some external
 1621 standard is not specified. This unspecified second is nominally equal to an International System
 1622 (SI) second in duration. Applications must be matched to a system that provides the particular
 1623 handling of external time in the way required by the application.

1624 IEEE Std 1003.1-2001/Cor 2-2004, item XBD/TC2/D6/12 is applied, making an editorial
 1625 correction to the paragraph commencing “How any changes to the value of seconds ...”.

1626 **A.4.16 Semaphore**

1627 There is no additional rationale provided for this section.

1628 **A.4.17 Thread-Safety**

1629 Where the interface of a function required by IEEE Std 1003.1-200x precludes thread-safety, an
 1630 alternate thread-safe form is provided. The names of these thread-safe forms are the same as the
 1631 non-thread-safe forms with the addition of the suffix “_r”. The suffix “_r” is historical, where
 1632 the ‘r’ stood for “reentrant”.

1633 In some cases, thread-safety is provided by restricting the arguments to an existing function.

1634 See also [Section B.2.9.1](#) (on page 163).

1635 **A.4.18 Tracing**

1636 Refer to [Section B.2.11](#) (on page 178).

1637 **A.4.19 Treatment of Error Conditions for Mathematical Functions**

1638 There is no additional rationale provided for this section.

1639 **A.4.20 Treatment of NaN Arguments for Mathematical Functions**

1640 There is no additional rationale provided for this section.

1641 **A.4.21 Utility**

1642 There is no additional rationale provided for this section.

1643 **A.4.22 Variable Assignment**

1644 There is no additional rationale provided for this section.

1645 A.5 File Format Notation

1646 The notation for spaces allows some flexibility for application output. Note that an empty
 1647 character position in *format* represents one or more <blank>s on the output (not *white space*,
 1648 which can include <newline>s). Therefore, another utility that reads that output as its input
 1649 must be prepared to parse the data using *scanf()*, *awk*, and so on. The 'Δ' character is used
 1650 when exactly one <space> is output.

1651 The treatment of integers and spaces is different from the *printf()* function in that they can be
 1652 surrounded with <blank>s. This was done so that, given a format such as:

```
1653 "%d\n" , <foo>
```

1654 the implementation could use a *printf()* call such as:

```
1655 printf("%6d\n" , foo);
```

1656 and still conform. This notation is thus somewhat like *scanf()* in addition to *printf()*.

1657 The *printf()* function was chosen as a model because most of the standard developers were
 1658 familiar with it. One difference from the C function *printf()* is that the l and h conversion
 1659 specifier characters are not used. As expressed by the Shell and Utilities volume of
 1660 IEEE Std 1003.1-200x, there is no differentiation between decimal values for type **int**, type **long**,
 1661 or type **short**. The conversion specifications %d or %i should be interpreted as an arbitrary
 1662 length sequence of digits. Also, no distinction is made between single precision and double
 1663 precision numbers (**float** or **double** in C). These are simply referred to as floating-point
 1664 numbers.

1665 Many of the output descriptions in the Shell and Utilities volume of IEEE Std 1003.1-200x use
 1666 the term "line", such as:

```
1667 "%s" , <input line>
```

1668 Since the definition of *line* includes the trailing <newline> already, there is no need to include a
 1669 '\n' in the format; a double <newline> would otherwise result.

1670 A.6 Character Set

1671 A.6.1 Portable Character Set

1672 The portable character set is listed in full so there is no dependency on the ISO/IEC 646:1991
 1673 standard (or historically ASCII) encoded character set, although the set is identical to the
 1674 characters defined in the International Reference version of the ISO/IEC 646:1991 standard.

1675 IEEE Std 1003.1-200x poses no requirement that multiple character sets or codesets be
 1676 supported, leaving this as a marketing differentiation for implementors. Although multiple
 1677 charmap files are supported, it is the responsibility of the implementation to provide the file(s);
 1678 if only one is provided, only that one will be accessible using the *localedef -f* option.

1679 The statement about invariance in codesets for the portable character set is worded to avoid
 1680 precluding implementations where multiple incompatible codesets are available (for instance,
 1681 ASCII and EBCDIC). The standard utilities cannot be expected to produce predictable results if
 1682 they access portable characters that vary on the same implementation.

1683 Not all character sets need include the portable character set, but each locale must include it. For
 1684 example, a Japanese-based locale might be supported by a mixture of character sets: JIS X 0201
 1685 Roman (a Japanese version of the ISO/IEC 646:1991 standard), JIS X 0208, and JIS X 0201

1686 Katakana. Not all of these character sets include the portable characters, but at least one does
1687 (JIS X 0201 Roman).

1688 A.6.2 Character Encoding

1689 Encoding mechanisms based on single shifts, such as the EUC encoding used in some Asian and
1690 other countries, can be supported via the current charmap mechanism. With single-shift
1691 encoding, each character is preceded by a shift code (SS2 or SS3). A complete EUC code,
1692 consisting of the portable character set (G0) and up to three additional character sets (G1, G2,
1693 G3), can be described using the current charmap mechanism; the encoding for each character in
1694 additional character sets G2 and G3 must then include their single-shift code. Other mechanisms
1695 to support locales based on encoding mechanisms such as locking shift are not addressed by this
1696 volume of IEEE Std 1003.1-200x.

1697 A.6.3 C Language Wide-Character Codes

1698 The standard does not specify how wide characters are encoded or provide a method for
1699 defining wide characters in a charmap. It specifies ways of translating between wide characters
1700 and multi-byte characters. The standard does not prevent an extension from providing a method
1701 to define wide characters.

1702 IEEE Std 1003.1-2001/Cor 2-2004, item XBD/TC2/D6/13 is applied, adding a statement that the
1703 standard has no means of defining a wide-character codeset.

1704 A.6.4 Character Set Description File

1705 IEEE PASC Interpretation 1003.2 #196 is applied, removing three lines of text dealing with
1706 ranges of symbolic names using position constant values which had been erroneously included
1707 in the final IEEE P1003.2b draft standard.

1708 IEEE Std 1003.1-2001/Cor 2-2004, item XBD/TC2/D6/14 is applied, correcting the example and
1709 adding a statement that the standard provides no means of defining a wide-character codeset.

1710 IEEE Std 1003.1-2001/Cor 2-2004, item XBD/TC2/D6/15 is applied, allowing the value zero for
1711 the width value of **WIDTH** and **WIDTH_DEFAULT**. This is required to cover some existing
1712 locales.

1713 A.6.4.1 State-Dependent Character Encodings

1714 A requirement was considered that would force utilities to eliminate any redundant locking
1715 shifts, but this was left as a quality of implementation issue.

1716 This change satisfies the following requirement from the ISO POSIX-2:1993 standard, Annex
1717 H.1:

1718 *The support of state-dependent (shift encoding) character sets should be addressed fully. See*
1719 *descriptions of these in the Base Definitions volume of IEEE Std 1003.1-200x, Section 6.2,*
1720 *Character Encoding. If such character encodings are supported, it is expected that this will impact*
1721 *the Base Definitions volume of IEEE Std 1003.1-200x, Section 6.2, Character Encoding, the Base*
1722 *Definitions volume of IEEE Std 1003.1-200x, Chapter 7, Locale, the Base Definitions volume of*
1723 *IEEE Std 1003.1-200x, Chapter 9, Regular Expressions , and the comm, cut, diff, grep, head, join,*
1724 *paste, and tail utilities.*

1725 The character set description file provides:

- 1726 • The capability to describe character set attributes (such as collation order or character
- 1727 classes) independent of character set encoding, and using only the characters in the
- 1728 portable character set. This makes it possible to create generic *localedef* source files for all
- 1729 codesets that share the portable character set (such as the ISO 8859 family or IBM Extended

1730 ASCII).

- 1731 • Standardized symbolic names for all characters in the portable character set, making it
- 1732 possible to refer to any such character regardless of encoding.

1733 Implementations are free to choose their own symbolic names, as long as the names identified
1734 by the Base Definitions volume of IEEE Std 1003.1-200x are also defined; this provides support
1735 for already existing “character names”.

1736 The names selected for the members of the portable character set follow the
1737 ISO/IEC 8859-1:1998 standard and the ISO/IEC 10646-1:2000 standard. However, several
1738 commonly used UNIX system names occur as synonyms in the list:

- 1739 • The historical UNIX system names are used for control characters.
- 1740 • The word “slash” is given in addition to “solidus”.
- 1741 • The word “backslash” is given in addition to “reverse-solidus”.
- 1742 • The word “hyphen” is given in addition to “hyphen-minus”.
- 1743 • The word “period” is given in addition to “full-stop”.
- 1744 • For digits, the word “digit” is eliminated.
- 1745 • For letters, the words “Latin Capital Letter” and “Latin Small Letter” are eliminated.
- 1746 • The words “left brace” and “right brace” are given in addition to “left-curly-bracket” and
1747 “right-curly-bracket”.
- 1748 • The names of the digits are preferred over the numbers to avoid possible confusion
1749 between ‘0’ and ‘O’, and between ‘1’ and ‘l’ (one and the letter ell).

1750 The names for the control characters in the Base Definitions volume of IEEE Std 1003.1-200x,
1751 Chapter 6, Character Set were taken from the ISO/IEC 4873:1991 standard.

1752 The charmap file was introduced to resolve problems with the portability of, especially, *localedef*
1753 sources. IEEE Std 1003.1-200x assumes that the portable character set is constant across all
1754 locales, but does not prohibit implementations from supporting two incompatible codings, such
1755 as both ASCII and EBCDIC. Such dual-support implementations should have all charmaps and
1756 *localedef* sources encoded using one portable character set, in effect cross-compiling for the other
1757 environment. Naturally, charmaps (and *localedef* sources) are only portable without
1758 transformation between systems using the same encodings for the portable character set. They
1759 can, however, be transformed between two sets using only a subset of the actual characters (the
1760 portable character set). However, the particular coded character set used for an application or an
1761 implementation does not necessarily imply different characteristics or collation; on the contrary,
1762 these attributes should in many cases be identical, regardless of codeset. The charmap provides
1763 the capability to define a common locale definition for multiple codesets (the same *localedef*
1764 source can be used for codesets with different extended characters; the ability in the charmap to
1765 define empty names allows for characters missing in certain codesets).

1766 The **<escape_char>** declaration was added at the request of the international community to ease
1767 the creation of portable charmap files on terminals not implementing the default backslash
1768 escape. The **<comment_char>** declaration was added at the request of the international
1769 community to eliminate the potential confusion between the number sign and the pound sign.

1770 The octal number notation with no leading zero required was selected to match those of *awk* and
1771 *tr* and is consistent with that used by *localedef*. To avoid confusion between an octal constant and
1772 the back-references used in *localedef* source, the octal, hexadecimal, and decimal constants must
1773 contain at least two digits. As single-digit constants are relatively rare, this should not impose
1774 any significant hardship. Provision is made for more digits to account for systems in which the

1775 byte size is larger than 8 bits. For example, a Unicode (ISO/IEC 10646-1:2000 standard) system
1776 that has defined 16-bit bytes may require six octal, four hexadecimal, and five decimal digits.

1777 The decimal notation is supported because some newer international standards define character
1778 values in decimal, rather than in the old column/row notation.

1779 The charmap identifies the coded character sets supported by an implementation. At least one
1780 charmap must be provided, but no implementation is required to provide more than one.
1781 Likewise, implementations can allow users to generate new charmaps (for instance, for a new
1782 version of the ISO 8859 family of coded character sets), but does not have to do so. If users are
1783 allowed to create new charmaps, the system documentation describes the rules that apply (for
1784 instance, “only coded character sets that are supersets of the ISO/IEC 646: 1991 standard IRV, no
1785 multi-byte characters”).

1786 This addition of the **WIDTH** specification satisfies the following requirement from the
1787 ISO POSIX-2: 1993 standard, Annex H.1:

1788 (9) *The definition of column position relies on the implementation’s knowledge of the integral*
1789 *width of the characters. The charmap or LC_CTYPE locale definitions should be enhanced to*
1790 *allow application specification of these widths.*

1791 The character “width” information was first considered for inclusion under *LC_CTYPE* but was
1792 moved because it is more closely associated with the information in the charmap than
1793 information in the locale source (cultural conventions information). Concerns were raised that
1794 formalizing this type of information is moving the locale source definition from the codeset-
1795 independent entity that it was designed to be to a repository of codeset-specific information. A
1796 similar issue occurred with the `<code_set_name>`, `<mb_cur_max>`, and `<mb_cur_min>`
1797 information, which was resolved to reside in the charmap definition.

1798 The width definition was added to the IEEE P1003.2b draft standard with the intent that the
1799 *wcswidth()* and/or *wcwidth()* functions (currently specified in the System Interfaces volume of
1800 IEEE Std 1003.1-200x) be the mechanism to retrieve the character width information.

1801 **A.7 Locale**

1802 **A.7.1 General**

1803 The description of locales is based on work performed in the UniForum Technical Committee,
1804 Subcommittee on Internationalization. Wherever appropriate, keywords are taken from the
1805 ISO C standard or the X/Open Portability Guide.

1806 The value used to specify a locale with environment variables is the name specified as the *name*
1807 operand to the *localedef* utility when the locale was created. This provides a verifiable method to
1808 create and invoke a locale.

1809 The “object” definitions need not be portable, as long as “source” definitions are. Strictly
1810 speaking, source definitions are portable only between implementations using the same
1811 character set(s). Such source definitions, if they use symbolic names only, easily can be ported
1812 between systems using different codesets, as long as the characters in the portable character set
1813 (see the Base Definitions volume of IEEE Std 1003.1-200x, Section 6.1, Portable Character Set)
1814 have common values between the codesets; this is frequently the case in historical
1815 implementations. Of source, this requires that the symbolic names used for characters outside
1816 the portable character set be identical between character sets. The definition of symbolic names
1817 for characters is outside the scope of IEEE Std 1003.1-200x, but is certainly within the scope of

1818 other standards organizations.

1819 Applications can select the desired locale by invoking the *setlocale()* function (or equivalent)
 1820 with the appropriate value. If the function is invoked with an empty string, the value of the
 1821 corresponding environment variable is used. If the environment variable is not set or is set to the
 1822 empty string, the implementation sets the appropriate environment as defined in the Base
 1823 Definitions volume of IEEE Std 1003.1-200x, Chapter 8, Environment Variables.

1824 A.7.2 POSIX Locale

1825 The POSIX locale is equal to the C locale. To avoid being classified as a C-language function, the
 1826 name has been changed to the POSIX locale; the environment variable value can be either
 1827 "POSIX" or, for historical reasons, "C".

1828 The POSIX definitions mirror the historical UNIX system behavior.

1829 The use of symbolic names for characters in the tables does not imply that the POSIX locale must
 1830 be described using symbolic character names, but merely that it may be advantageous to do so.

1831 A.7.3 Locale Definition

1832 The decision to separate the file format from the *localedef* utility description was only partially
 1833 editorial. Implementations may provide other interfaces than *localedef*. Requirements on "the
 1834 utility", mostly concerning error messages, are described in this way because they are meant to
 1835 affect the other interfaces implementations may provide as well as *localedef*.

1836 The text about POSIX2_LOCALEDEF does not mean that internationalization is optional; only
 1837 that the functionality of the *localedef* utility is. REs, for instance, must still be able to recognize,
 1838 for example, character class expressions such as "[[:alpha:]]". A possible analogy is with
 1839 an applications development environment; while all conforming implementations must be
 1840 capable of executing applications, not all need to have the development environment installed.
 1841 The assumption is that the capability to modify the behavior of utilities (and applications) via
 1842 locale settings must be supported. If the *localedef* utility is not present, then the only choice is to
 1843 select an existing (presumably implementation-documented) locale. An implementation could,
 1844 for example, choose to support only the POSIX locale, which would in effect limit the amount of
 1845 changes from historical implementations quite drastically. The *localedef* utility is still required,
 1846 but would always terminate with an exit code indicating that no locale could be created.
 1847 Supported locales must be documented using the syntax defined in this chapter. (This ensures
 1848 that users can accurately determine what capabilities are provided. If the implementation
 1849 decides to provide additional capabilities to the ones in this chapter, that is already provided
 1850 for.)

1851 If the option is present (that is, locales can be created), then the *localedef* utility must be capable
 1852 of creating locales based on the syntax and rules defined in this chapter. This does not mean that
 1853 the implementation cannot also provide alternate means for creating locales.

1854 The octal, decimal, and hexadecimal notations are the same employed by the charmap facility
 1855 (see the Base Definitions volume of IEEE Std 1003.1-200x, Section 6.4, Character Set Description
 1856 File). To avoid confusion between an octal constant and a back-reference, the octal, hexadecimal,
 1857 and decimal constants must contain at least two digits. As single-digit constants are relatively
 1858 rare, this should not impose any significant hardship. Provision is made for more digits to
 1859 account for systems in which the byte size is larger than 8 bits. For example, a Unicode (see the
 1860 ISO/IEC 10646-1:2000 standard) system that has defined 16-bit bytes may require six octal, four
 1861 hexadecimal, and five decimal digits. As with the charmap file, multi-byte characters are
 1862 described in the locale definition file using "big-endian" notation for reasons of portability.
 1863 There is no requirement that the internal representation in the computer memory be in this same
 1864 order.

1865 One of the guidelines used for the development of this volume of IEEE Std 1003.1-200x is that
 1866 characters outside the invariant part of the ISO/IEC 646:1991 standard should not be used in
 1867 portable specifications. The backslash character is not in the invariant part; the number sign is,
 1868 but with multiple representations: as a number sign, and as a pound sign. As far as general
 1869 usage of these symbols, they are covered by the “grandfather clause”, but for newly defined
 1870 interfaces, the WG15 POSIX working group has requested that POSIX provide alternate
 1871 representations. Consequently, while the default escape character remains the backslash and the
 1872 default comment character is the number sign, implementations are required to recognize
 1873 alternative representations, identified in the applicable source file via the `<escape_char>` and
 1874 `<comment_char>` keywords.

1875 A.7.3.1 *LC_CTYPE*

1876 The *LC_CTYPE* category is primarily used to define the encoding-independent aspects of a
 1877 character set, such as character classification. In addition, certain encoding-dependent
 1878 characteristics are also defined for an application via the *LC_CTYPE* category.
 1879 IEEE Std 1003.1-200x does not mandate that the encoding used in the locale is the same as the
 1880 one used by the application because an implementation may decide that it is advantageous to
 1881 define locales in a system-wide encoding rather than having multiple, logically identical locales
 1882 in different encodings, and to convert from the application encoding to the system-wide
 1883 encoding on usage. Other implementations could require encoding-dependent locales.

1884 In either case, the *LC_CTYPE* attributes that are directly dependent on the encoding, such as
 1885 `<mb_cur_max>` and the display width of characters, are not user-specifiable in a locale source
 1886 and are consequently not defined as keywords.

1887 Implementations may define additional keywords or extend the *LC_CTYPE* mechanism to allow
 1888 application-defined keywords.

1889 The text “The ellipsis specification shall only be valid within a single encoded character set” is
 1890 present because it is possible to have a locale supported by multiple character encodings, as
 1891 explained in the rationale for the Base Definitions volume of IEEE Std 1003.1-200x, Section 6.1,
 1892 Portable Character Set. An example given there is of a possible Japanese-based locale supported
 1893 by a mixture of the character sets JIS X 0201 Roman, JIS X 0208, and JIS X 0201 Katakana.
 1894 Attempting to express a range of characters across these sets is not logical and the
 1895 implementation is free to reject such attempts.

1896 As the *LC_CTYPE* character classes are based on the ISO C standard character class definition,
 1897 the category does not support multi-character elements. For instance, the German character
 1898 `<sharp-s>` is traditionally classified as a lowercase letter. There is no corresponding uppercase
 1899 letter; in proper capitalization of German text, the `<sharp-s>` will be replaced by “SS”; that is, by
 1900 two characters. This kind of conversion is outside the scope of the **toupper** and **tolower**
 1901 keywords.

1902 Where IEEE Std 1003.1-200x specifies that only certain characters can be specified, as for the
 1903 keywords **digit** and **xdigit**, the specified characters must be from the portable character set, as
 1904 shown. As an example, only the Arabic digits 0 through 9 are acceptable as digits.

1905 The character classes **digit**, **xdigit**, **lower**, **upper**, and **space** have a set of automatically included
 1906 characters. These only need to be specified if the character values (that is, encoding) differs from
 1907 the implementation default values. It is not possible to define a locale without these
 1908 automatically included characters unless some implementation extension is used to prevent
 1909 their inclusion. Such a definition would not be a proper superset of the C locale, and thus, it
 1910 might not be possible for the standard utilities to be implemented as programs conforming to
 1911 the ISO C standard.

1912 The definition of character class **digit** requires that only ten characters—the ones defining

1913 digits—can be specified; alternate digits (for example, Hindi or Kanji) cannot be specified here.
 1914 However, the encoding may vary if an implementation supports more than one encoding.

1915 The definition of character class **xdigit** requires that the characters included in character class
 1916 **digit** are included here also and allows for different symbols for the hexadecimal digits 10
 1917 through 15.

1918 The inclusion of the **charclass** keyword satisfies the following requirement from the
 1919 ISO POSIX-2: 1993 standard, Annex H.1:

1920 (3) *The LC_CTYPE (2.5.2.1) locale definition should be enhanced to allow user-specified additional*
 1921 *character classes, similar in concept to the ISO C standard Multibyte Support Extension (MSE)*
 1922 *iswctype() function.*

1923 This keyword was previously included in The Open Group specifications and is now mandated
 1924 in the Shell and Utilities volume of IEEE Std 1003.1-200x.

1925 The symbolic constant {CHARCLASS_NAME_MAX} was also adopted from The Open Group
 1926 specifications. Applications portability is enhanced by the use of symbolic constants.

1927 A.7.3.2 LC_COLLATE

1928 The rules governing collation depend to some extent on the use. At least five different levels of
 1929 increasingly complex collation rules can be distinguished:

- 1930 1. *Byte/machine code order*: This is the historical collation order in the UNIX system and many
 1931 proprietary operating systems. Collation is here performed character by character,
 1932 without any regard to context. The primary virtue is that it usually is quite fast and also
 1933 completely deterministic; it works well when the native machine collation sequence
 1934 matches the user expectations.
- 1935 2. *Character order*: On this level, collation is also performed character by character, without
 1936 regard to context. The order between characters is, however, not determined by the code
 1937 values, but on the expectations by the user of the “correct” order between characters. In
 1938 addition, such a (simple) collation order can specify that certain characters collate equally
 1939 (for example, uppercase and lowercase letters).
- 1940 3. *String ordering*: On this level, entire strings are compared based on relatively
 1941 straightforward rules. Several “passes” may be required to determine the order between
 1942 two strings. Characters may be ignored in some passes, but not in others; the strings may
 1943 be compared in different directions; and simple string substitutions may be performed
 1944 before strings are compared. This level is best described as “dictionary” ordering; it is
 1945 based on the spelling, not the pronunciation, or meaning, of the words.
- 1946 4. *Text search ordering*: This is a further refinement of the previous level, best described as
 1947 “telephone book ordering”; some common homonyms (words spelled differently but
 1948 with the same pronunciation) are collated together; numbers are collated as if they were
 1949 spelled out, and so on.
- 1950 5. *Semantic-level ordering*: Words and strings are collated based on their meaning; entire
 1951 words (such as “the”) are eliminated; the ordering is not deterministic. This usually
 1952 requires special software and is highly dependent on the intended use.

1953 While the historical collation order formally is at level 1, for the English language it corresponds
 1954 roughly to elements at level 2. The user expects to see the output from the *ls* utility sorted very
 1955 much as it would be in a dictionary. While telephone book ordering would be an optimal goal
 1956 for standard collation, this was ruled out as the order would be language-dependent.
 1957 Furthermore, a requirement was that the order must be determined solely from the text string

1958 and the collation rules; no external information (for example, “pronunciation dictionaries”)
 1959 could be required.

1960 As a result, the goal for the collation support is at level 3. This also matches the requirements for
 1961 the Canadian collation order, as well as other, known collation requirements for alphabetic
 1962 scripts. It specifically rules out collation based on pronunciation rules or based on semantic
 1963 analysis of the text.

1964 The syntax for the *LC_COLLATE* category source meets the requirements for level 3 and has
 1965 been verified to produce the correct result with examples based on French, Canadian, and
 1966 Danish collation order. Because it supports multi-character collating elements, it is also capable
 1967 of supporting collation in codesets where a character is expressed using non-spacing characters
 1968 followed by the base character (such as the ISO/IEC 6937: 2001 standard).

1969 The directives that can be specified in an operand to the **order_start** keyword are based on the
 1970 requirements specified in several proposed standards and in customary use. The following is a
 1971 rephrasing of rules defined for “lexical ordering in English and French” by the Canadian
 1972 Standards Association (the text in square brackets is rephrased):

- 1973 • Once special characters [punctuation] have been removed from original strings, the
 1974 ordering is determined by scanning forwards (left to right) [disregarding case and
 1975 diacriticals].
- 1976 • In case of equivalence, special characters are once again removed from original strings and
 1977 the ordering is determined by scanning backwards (starting from the rightmost character
 1978 of the string and back), character by character [disregarding case but considering
 1979 diacriticals].
- 1980 • In case of repeated equivalence, special characters are removed again from original strings
 1981 and the ordering is determined by scanning forwards, character by character [considering
 1982 both case and diacriticals].
- 1983 • If there is still an ordering equivalence after the first three rules have been applied, then
 1984 only special characters and the position they occupy in the string are considered to
 1985 determine ordering. The string that has a special character in the lowest position comes
 1986 first. If two strings have a special character in the same position, the character [with the
 1987 lowest collation value] comes first. In case of equality, the other special characters are
 1988 considered until there is a difference or until all special characters have been exhausted.

1989 It is estimated that this part of IEEE Std 1003.1-200x covers the requirements for all European
 1990 languages, and no particular problems are anticipated with Slavic or Middle East character sets.

1991 The Far East (particularly Japanese/Chinese) collations are often based on contextual
 1992 information and pronunciation rules (the same ideogram can have different meanings and
 1993 different pronunciations). Such collation, in general, falls outside the desired goal of
 1994 IEEE Std 1003.1-200x. There are, however, several other collation rules (stroke/radical or “most
 1995 common pronunciation”) that can be supported with the mechanism described here.

1996 The character order is defined by the order in which characters and elements are specified
 1997 between the **order_start** and **order_end** keywords. Weights assigned to the characters and
 1998 elements define the collation sequence; in the absence of weights, the character order is also the
 1999 collation sequence.

2000 The **position** keyword provides the capability to consider, in a compare, the relative position of
 2001 characters not subject to **IGNORE**. As an example, consider the two strings "o-ring" and
 2002 "or-ing". Assuming the hyphen is subject to **IGNORE** on the first pass, the two strings
 2003 compare equal, and the position of the hyphen is immaterial. On second pass, all characters
 2004 except the hyphen are subject to **IGNORE**, and in the normal case the two strings would again

2005 compare equal. By taking position into account, the first collates before the second.

2006 A.7.3.3 LC_MONETARY

2007 The currency symbol does not appear in *LC_MONETARY* because it is not defined in the C
2008 locale of the ISO C standard.

2009 The ISO C standard limits the size of decimal points and thousands delimiters to single-byte
2010 values. In locales based on multi-byte coded character sets, this cannot be enforced;
2011 IEEE Std 1003.1-200x does not prohibit such characters, but makes the behavior unspecified (in
2012 the text “In contexts where other standards ...”).

2013 The grouping specification is based on, but not identical to, the ISO C standard. The -1 indicates
2014 that no further grouping is performed; the equivalent of {CHAR_MAX} in the ISO C standard.

2015 The text “the value is not available in the locale” is taken from the ISO C standard and is used
2016 instead of the “unspecified” text in early proposals. There is no implication that omitting these
2017 keywords or assigning them values of " " or -1 produces unspecified results; such omissions or
2018 assignments eliminate the effects described for the keyword or produce zero-length strings, as
2019 appropriate.

2020 The locale definition is an extension of the ISO C standard *localeconv()* specification. In
2021 particular, rules on how **currency_symbol** is treated are extended to also cover **int_curr_symbol**,
2022 and **p_sep_by_space** and **n_sep_by_space** have been augmented with the value 2, which places
2023 a <space> between the sign and the symbol. This has been updated to match the
2024 ISO/IEC 9899:1999 standard requirements and is an incompatible change from UNIX 98 and the
2025 ISO POSIX-2 standard and the ISO POSIX-1:1996 standard requirements. The following table
2026 shows the result of various combinations:

		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\$+

2039 The following is an example of the interpretation of the **mon_grouping** keyword. Assuming that
2040 the value to be formatted is 123 456 789 and the **mon_thousands_sep** is ' ', then the following
2041 table shows the result. The third column shows the equivalent string in the ISO C standard that
2042 would be used by the *localeconv()* function to accommodate this grouping.

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"

2049 In these examples, the octal value of {CHAR_MAX} is 177.

2050 IEEE Std 1003.1-2001/Cor 1-2002, item XBD/TC1/D6/6 adds a correction that permits the Euro

2051 currency symbol and addresses extensibility. The correction is stated using the term “should”
 2052 intentionally, in order to make this a recommendation rather than a restriction on
 2053 implementations. This allows for flexibility in implementations on how they handle future
 2054 currency symbol additions.

2055 IEEE Std 1003.1-2001/Cor 1-2002, tem XBD/TC1/D6/5 is applied, adding the **int_[np]_*** values
 2056 to the POSIX locale definition of *LC_MONETARY*.

2057 IEEE Std 1003.1-2001/Cor 2-2004, item XBD/TC2/D6/16 is applied, updating the descriptions
 2058 of **p_sep_by_space**, **n_sep_by_space**, **int_p_sep_by_space**, and **int_n_sep_by_space** to match
 2059 the description of these keywords in the ISO C standard and the System Interfaces volume of
 2060 IEEE Std 1003.1-200x, *localeconv()*.

2061 A.7.3.4 *LC_NUMERIC*

2062 See the rationale for *LC_MONETARY* for a description of the behavior of grouping.

2063 A.7.3.5 *LC_TIME*

2064 Although certain of the conversion specifications in the POSIX locale (such as the name of the
 2065 month) are shown with initial capital letters, this need not be the case in other locales. Programs
 2066 using these conversion specifications may need to adjust the capitalization if the output is going
 2067 to be used at the beginning of a sentence.

2068 The *LC_TIME* descriptions of **abday**, **day**, **mon**, and **abmon** imply a Gregorian style calendar
 2069 (7-day weeks, 12-month years, leap years, and so on). Formatting time strings for other types of
 2070 calendars is outside the scope of IEEE Std 1003.1-200x.

2071 While the ISO 8601:2000 standard numbers the weekdays starting with Monday, historical
 2072 practice is to use the Sunday as the first day. Rather than change the order and introduce
 2073 potential confusion, the days must be specified beginning with Sunday; previous references to
 2074 “first day” have been removed. Note also that the Shell and Utilities volume of
 2075 IEEE Std 1003.1-200x *date* utility supports numbering compliant with the ISO 8601:2000
 2076 standard.

2077 As specified under *date* in the Shell and Utilities volume of IEEE Std 1003.1-200x and *strftime()*
 2078 in the System Interfaces volume of IEEE Std 1003.1-200x, the conversion specifications
 2079 corresponding to the optional keywords consist of a modifier followed by a traditional
 2080 conversion specification (for instance, %Ex). If the optional keywords are not supported by the
 2081 implementation or are unspecified for the current locale, these modified conversion
 2082 specifications are treated as the traditional conversion specifications. For example, assume the
 2083 following keywords:

```
2084 alt_digits "0th"; "1st"; "2nd"; "3rd"; "4th"; "5th"; \  

  2085           "6th"; "7th"; "8th"; "9th"; "10th"
```

```
2086 d_fmt     "The %Od day of %B in %Y"
```

2087 On July 4th 1776, the %x conversion specifications would result in "The 4th day of July
 2088 in 1776", while on July 14th 1789 it would result in "The 14 day of July in 1789". It
 2089 can be noted that the above example is for illustrative purposes only; the %O modifier is
 2090 primarily intended to provide for Kanji or Hindi digits in *date* formats.

2091 The following is an example for Japan that supports the current plus last three Emperors and
 2092 reverts to Western style numbering for years prior to the Meiji era. The example also allows for
 2093 the custom of using a special name for the first year of an era instead of using 1. (The examples
 2094 substitute romaji where kanji should be used.)

```
2095 era_d_fmt "%EY%mgatsu%dnichi (%a)"
```

```

2096 era    "+:2:1990/01/01:+*:Heisei:%EC%Eynen";\
2097      "+:1:1989/01/08:1989/12/31:Heisei:%ECgannen";\
2098      "+:2:1927/01/01:1989/01/07:Shouwa:%EC%Eynen";\
2099      "+:1:1926/12/25:1926/12/31:Shouwa:%ECgannen";\
2100      "+:2:1913/01/01:1926/12/24:Taishou:%EC%Eynen";\
2101      "+:1:1912/07/30:1912/12/31:Taishou:%ECgannen";\
2102      "+:2:1869/01/01:1912/07/29:Meiji:%EC%Eynen";\
2103      "+:1:1868/09/08:1868/12/31:Meiji:%ECgannen";\
2104      "-:1868:1868/09/07:-*::%Ey"

```

2105 Assuming that the current date is September 21, 1991, a request to *date* or *strftime()* would yield
 2106 the following results:

```

2107 %Ec - Heisei3nen9gatsu2lnichi (Sat) 14:39:26
2108 %EC - Heisei
2109 %Ex - Heisei3nen9gatsu2lnichi (Sat)
2110 %Ey - 3
2111 %EY - Heisei3nen

```

2112 Example era definitions for the Republic of China:

```

2113 era    "+:2:1913/01/01:+*:ChungHwaMingGuo:%EC%EyNen";\
2114      "+:1:1912/1/1:1912/12/31:ChungHwaMingGuo:%ECYuenNen";\
2115      "+:1:1911/12/31:-*:MingChien:%EC%EyNen"

```

2116 Example definitions for the Christian Era:

```

2117 era    "+:1:0001/01/01:+*:AD:%EC %Ey";\
2118      "+:1:-0001/12/31:-*:BC:%Ey %EC"

```

2119 A.7.3.6 LC_MESSAGES

2120 The **yesstr** and **nostr** locale keywords and the YESSTR and NOSTR *langinfo* items were formerly
 2121 used to match user affirmative and negative responses. In IEEE Std 1003.1-200x, the **yesexpr**,
 2122 **noexpr**, YESEXPR, and NOEXPR extended regular expressions have replaced them.
 2123 Applications should use the general locale-based messaging facilities to issue prompting
 2124 messages which include sample desired responses.

2125 A.7.4 Locale Definition Grammar

2126 There is no additional rationale provided for this section.

2127 A.7.4.1 Locale Lexical Conventions

2128 There is no additional rationale provided for this section.

2129 A.7.4.2 Locale Grammar

2130 There is no additional rationale provided for this section.

2131 A.7.5 Locale Definition Example

2132 The following is an example of a locale definition file that could be used as input to the *localedef*
 2133 utility. It assumes that the utility is executed with the *-f* option, naming a charmap file with (at
 2134 least) the following content:

```

2135 CHARMAP
2136 <space>      \x20
2137 <dollar>     \x24
2138 <A>          \101
2139 <a>          \141

```



```

2140 <A-acute>    \346
2141 <a-acute>    \365
2142 <A-grave>    \300
2143 <a-grave>    \366
2144 <b>          \142
2145 <C>          \103
2146 <c>          \143
2147 <c-cedilla>  \347
2148 <d>          \x64
2149 <H>          \110
2150 <h>          \150
2151 <eszet>     \xb7
2152 <s>          \x73
2153 <z>          \x7a
2154 END CHARMAP

```

2155 It should not be taken as complete or to represent any actual locale, but only to illustrate the
 2156 syntax.

```

2157 #
2158 LC_CTYPE
2159 lower <a>;<b>;<c>;<c-cedilla>;<d>;...;<z>
2160 upper A;B;C;Ç;...;Z
2161 space \x20;\x09;\x0a;\x0b;\x0c;\x0d
2162 blank \040;\011
2163 toupper (<a>,<A>);(b,B);(c,C);(ç,Ç);(d,D);(z,Z)
2164 END LC_CTYPE
2165 #
2166 LC_COLLATE
2167 #
2168 # The following example of collation is based on
2169 # Canadian standard Z243.4.1-1998, "Canadian Alphanumeric
2170 # Ordering Standard for Character Sets of CSA Z234.4 Standard".
2171 # (Other parts of this example locale definition file do not
2172 # purport to relate to Canada, or to any other real culture.)
2173 # The proposed standard defines a 4-weight collation, such that
2174 # in the first pass, characters are compared without regard to
2175 # case or accents; in the second pass, backwards-compare without
2176 # regard to case; in the third pass, forwards-compare without
2177 # regard to diacriticals. In the 3 first passes, non-alphabetic
2178 # characters are ignored; in the fourth pass, only special
2179 # characters are considered, such that "The string that has a
2180 # special character in the lowest position comes first. If two
2181 # strings have a special character in the same position, the
2182 # collation value of the special character determines ordering.
2183 #
2184 # Only a subset of the character set is used here; mostly to
2185 # illustrate the set-up.
2186 #
2187 collating-symbol <NULL>
2188 collating-symbol <LOW_VALUE>
2189 collating-symbol <LOWER-CASE>
2190 collating-symbol <SUBSCRIPT-LOWER>
2191 collating-symbol <SUPERSCRIPT-LOWER>

```



```

2192 collating-symbol <UPPER-CASE>
2193 collating-symbol <NO-ACCENT>
2194 collating-symbol <PECULIAR>
2195 collating-symbol <LIGATURE>
2196 collating-symbol <ACUTE>
2197 collating-symbol <GRAVE>
2198 # Further collating-symbols follow.
2199 #
2200 # Properly, the standard does not include any multi-character
2201 # collating elements; the one below is added for completeness.
2202 #
2203 collating_element <ch> from "<c><h>"
2204 collating_element <CH> from "<C><H>"
2205 collating_element <Ch> from "<C><h>"
2206 #
2207 order_start forward;backward;forward;forward,position
2208 #
2209 # Collating symbols are specified first in the sequence to allocate
2210 # basic collation values to them, lower than that of any character.
2211 <NULL>
2212 <LOW_VALUE>
2213 <LOWER-CASE>
2214 <SUBSCRIPT-LOWER>
2215 <SUPERSCRIPT-LOWER>
2216 <UPPER-CASE>
2217 <NO-ACCENT>
2218 <PECULIAR>
2219 <LIGATURE>
2220 <ACUTE>
2221 <GRAVE>
2222 <RING-ABOVE>
2223 <DIAERESIS>
2224 <TILDE>
2225 # Further collating symbols are given a basic collating value here.
2226 #
2227 # Here follow special characters.
2228 <space> IGNORE;IGNORE;IGNORE;<space>
2229 # Other special characters follow here.
2230 #
2231 # Here follow the regular characters.
2232 <a> <a>;<NO-ACCENT>;<LOWER-CASE>;IGNORE
2233 <A> <a>;<NO-ACCENT>;<UPPER-CASE>;IGNORE
2234 <a-acute> <a>;<ACUTE>;<LOWER-CASE>;IGNORE
2235 <A-acute> <a>;<ACUTE>;<UPPER-CASE>;IGNORE
2236 <a-grave> <a>;<GRAVE>;<LOWER-CASE>;IGNORE
2237 <A-grave> <a>;<GRAVE>;<UPPER-CASE>;IGNORE
2238 <ae> " <a><e> " ; " <LIGATURE><LIGATURE> " ; \
2239 " <LOWER-CASE><LOWER-CASE> " ; IGNORE
2240 <AE> " <a><e> " ; " <LIGATURE><LIGATURE> " ; \
2241 " <UPPER-CASE><UPPER-CASE> " ; IGNORE
2242 <b> <b>;<NO-ACCENT>;<LOWER-CASE>;IGNORE
2243 <B> <b>;<NO-ACCENT>;<UPPER-CASE>;IGNORE
2244 <c> <c>;<NO-ACCENT>;<LOWER-CASE>;IGNORE

```

```

2245 <C>          <c>; <NO-ACCENT>; <UPPER-CASE>; IGNORE
2246 <ch>         <ch>; <NO-ACCENT>; <LOWER-CASE>; IGNORE
2247 <Ch>         <ch>; <NO-ACCENT>; <PECULIAR>; IGNORE
2248 <CH>         <ch>; <NO-ACCENT>; <UPPER-CASE>; IGNORE
2249 #
2250 # As an example, the strings "Bach" and "bach" could be encoded (for
2251 # compare purposes) as:
2252 # "Bach"  <b>; <a>; <ch>; <LOW_VALUE>; <NO-ACCENT>; <NO-ACCENT>; \
2253 #         <NO-ACCENT>; <LOW_VALUE>; <UPPER-CASE>; <LOWER-CASE>; \
2254 #         <LOWER-CASE>; <NULL>
2255 # "bach"  <b>; <a>; <ch>; <LOW_VALUE>; <NO-ACCENT>; <NO-ACCENT>; \
2256 #         <NO-ACCENT>; <LOW_VALUE>; <LOWER-CASE>; <LOWER-CASE>; \
2257 #         <LOWER-CASE>; <NULL>
2258 #
2259 # The two strings are equal in pass 1 and 2, but differ in pass 3.
2260 #
2261 # Further characters follow.
2262 #
2263 UNDEFINED    IGNORE; IGNORE; IGNORE; IGNORE
2264 #
2265 order_end
2266 #
2267 END LC_COLLATE
2268 #
2269 LC_MONETARY
2270 int_curr_symbol    "USD "
2271 currency_symbol   "$"
2272 mon_decimal_point  "."
2273 mon_grouping      3;0
2274 positive_sign     ""
2275 negative_sign     "- "
2276 p_cs_precedes     1
2277 n_sign_posn       0
2278 END LC_MONETARY
2279 #
2280 LC_NUMERIC
2281 copy "US_en.ASCII"
2282 END LC_NUMERIC
2283 #
2284 LC_TIME
2285 abday  "Sun"; "Mon"; "Tue"; "Wed"; "Thu"; "Fri"; "Sat"
2286 #
2287 day    "Sunday"; "Monday"; "Tuesday"; "Wednesday"; \
2288        "Thursday"; "Friday"; "Saturday"
2289 #
2290 abmon  "Jan"; "Feb"; "Mar"; "Apr"; "May"; "Jun"; \
2291        "Jul"; "Aug"; "Sep"; "Oct"; "Nov"; "Dec"
2292 #
2293 mon    "January"; "February"; "March"; "April"; \
2294        "May"; "June"; "July"; "August"; "September"; \
2295        "October"; "November"; "December"
2296 #
2297 d_t_fmt "%a %b %d %T %Z %Y\n"

```

```

2298     END LC_TIME
2299     #
2300     LC_MESSAGES
2301     yesexpr  "^[yY][[:alpha:]]*" | (OK) "
2302     #
2303     noexpr   "^[nN][[:alpha:]]*"
2304     END LC_MESSAGES

```

2305 A.8 Environment Variables

2306 A.8.1 Environment Variable Definition

2307 The variable *environ* is not intended to be declared in any header, but rather to be declared by
 2308 the user for accessing the array of strings that is the environment. This is the traditional usage of
 2309 the symbol. Putting it into a header could break some programs that use the symbol for their
 2310 own purposes.

2311 The decision to restrict conforming systems to the use of digits, uppercase letters, and
 2312 underscores for environment variable names allows applications to use lowercase letters in their
 2313 environment variable names without conflicting with any conforming system.

2314 In addition to the obvious conflict with the shell syntax for positional parameter substitution,
 2315 some historical applications (including some shells) exclude names with leading digits from the
 2316 environment.

2317 A.8.2 Internationalization Variables

2318 Utilities conforming to the Shell and Utilities volume of IEEE Std 1003.1-200x and written in
 2319 standard C can access the locale variables by issuing the following call:

```
2320 setlocale(LC_ALL, "")
```

2321 If this were omitted, the ISO C standard specifies that the C locale would be used.

2322 The DESCRIPTION of *setlocale()* requires that when setting all categories of a locale, if the value
 2323 of any of the environment variable searches yields a locale that is not supported (and non-null),
 2324 the *setlocale()* function returns a null pointer and the locale of the process is unchanged.

2325 For the standard utilities, if any of the environment variables are invalid, it makes sense to
 2326 default to an implementation-defined, consistent locale environment. It is more confusing for a
 2327 user to have partial settings occur in case of a mistake. All utilities would then behave in one
 2328 language/cultural environment. Furthermore, it provides a way of forcing the whole
 2329 environment to be the implementation-defined default. Disastrous results could occur if a
 2330 pipeline of utilities partially uses the environment variables in different ways. In this case, it
 2331 would be appropriate for utilities that use *LANG* and related variables to exit with an error if
 2332 any of the variables are invalid. For example, users typing individual commands at a terminal
 2333 might want *date* to work if *LC_MONETARY* is invalid as long as *LC_TIME* is valid. Since these
 2334 are conflicting reasonable alternatives, IEEE Std 1003.1-200x leaves the results unspecified if the
 2335 locale environment variables would not produce a complete locale matching the specification of
 2336 the user.

2337 The locale settings of individual categories cannot be truly independent and still guarantee
 2338 correct results. For example, when collating two strings, characters must first be extracted from
 2339 each string (governed by *LC_CTYPE*) before being mapped to collating elements (governed by
 2340 *LC_COLLATE*) for comparison. That is, if *LC_CTYPE* is causing parsing according to the rules of

2341 a large, multi-byte code set (potentially returning 20 000 or more distinct character codeset
 2342 values), but *LC_COLLATE* is set to handle only an 8-bit codeset with 256 distinct characters,
 2343 meaningful results are obviously impossible.

2344 The *LC_MESSAGES* variable affects the language of messages generated by the standard
 2345 utilities.

2346 The description of the environment variable names starting with the characters “LC_”
 2347 acknowledges the fact that the interfaces presented may be extended as new international
 2348 functionality is required. In the ISO C standard, names preceded by “LC_” are reserved in the
 2349 name space for future categories.

2350 To avoid name clashes, new categories and environment variables are divided into two
 2351 classifications: “implementation-independent” and “implementation-defined”.

2352 Implementation-independent names will have the following format:

2353 *LC_NAME*

2354 where *NAME* is the name of the new category and environment variable. Capital letters must be
 2355 used for implementation-independent names.

2356 Implementation-defined names must be in lowercase letters, as below:

2357 *LC_name*

2358 **A.8.3 Other Environment Variables**

2359 **COLUMNS, LINES**

2360 The default values for the number of column positions, *COLUMNS*, and screen height, *LINES*,
 2361 are unspecified because historical implementations use different methods to determine values
 2362 corresponding to the size of the screen in which the utility is run. This size is typically known to
 2363 the implementation through the value of *TERM*, or by more elaborate methods such as
 2364 extensions to the *stty* utility or knowledge of how the user is dynamically resizing windows on a
 2365 bit-mapped display terminal. Users should not need to set these variables in the environment
 2366 unless there is a specific reason to override the default behavior of the implementation, such as
 2367 to display data in an area arbitrarily smaller than the terminal or window. Values for these
 2368 variables that are not decimal integers greater than zero are implicitly undefined values; it is
 2369 unnecessary to enumerate all of the possible values outside of the acceptable set.

2370 **LOGNAME**

2371 In most implementations, the value of such a variable is easily forged, so security-critical
 2372 applications should rely on other means of determining user identity. *LOGNAME* is required to
 2373 be constructed from the portable filename character set for reasons of interchange. No diagnostic
 2374 condition is specified for violating this rule, and no requirement for enforcement exists. The
 2375 intent of the requirement is that if extended characters are used, the “guarantee” of portability
 2376 implied by a standard is void.

2377 **PATH**

2378 Many historical implementations of the Bourne shell do not interpret a trailing colon to
 2379 represent the current working directory and are thus non-conforming. The C Shell and the
 2380 KornShell conform to IEEE Std 1003.1-200x on this point. The usual name of dot may also be
 2381 used to refer to the current working directory.

2382 Many implementations historically have used a default value of */bin* and */usr/bin* for the *PATH*
 2383 variable. IEEE Std 1003.1-200x does not mandate this default path be identical to that retrieved
 2384 from *getconf _CS_PATH* because it is likely that the standardized utilities may be provided in

2385 another directory separate from the directories used by some historical applications.

2386 SHELL

2387 The *SHELL* variable names the preferred shell of the user; it is a guide to applications. There is
 2388 no direct requirement that that shell conform to IEEE Std 1003.1-200x; that decision should rest
 2389 with the user. It is the intention of the standard developers that alternative shells be permitted, if
 2390 the user chooses to develop or acquire one. An operating system that builds its shell into the
 2391 “kernel” in such a manner that alternative shells would be impossible does not conform to the
 2392 spirit of IEEE Std 1003.1-200x.

2393 TZ

2394 The quoted form of the timezone variable allows timezone names of the form UTC+1 (or any
 2395 name that contains the character plus ('+'), the character minus ('-'), or digits), which may be
 2396 appropriate for countries that do not have an official timezone name. It would be coded as
 2397 <UTC+1>+1<UTC+2>, which would cause *std* to have a value of UTC+1 and *dst* a value of
 2398 UTC+2, each with a length of 5 characters. This does not appear to conflict with any existing
 2399 usage. The characters '<' and '>' were chosen for quoting because they are easier to parse
 2400 visually than a quoting character that does not provide some sense of bracketing (and in a string
 2401 like this, such bracketing is helpful). They were also chosen because they do not need special
 2402 treatment when assigning to the *TZ* variable. Users are often confused by embedding quotes in a
 2403 string. Because '<' and '>' are meaningful to the shell, the whole string would have to be
 2404 quoted, but that is easily explained. (Parentheses would have presented the same problems.)
 2405 Although the '>' symbol could have been permitted in the string by either escaping it or
 2406 doubling it, it seemed of little value to require that. This could be provided as an extension if
 2407 there was a need. Timezone names of this new form lead to a requirement that the value of
 2408 {_POSIX_TZNAME_MAX} change from 3 to 6.

2409 Since the *TZ* environment variable is usually inherited by all applications started by a user after
 2410 the value of the *TZ* environment variable is changed and since many applications run using the
 2411 C or POSIX locale, using characters that are not in the portable character set in the *std* and *dst*
 2412 fields could cause unexpected results.

2413 The format of the *TZ* environment variable is changed in Issue 6 to allow for the quoted form, as
 2414 defined in previous versions of the ISO POSIX-1 standard.

2415 IEEE Std 1003.1-2001/Cor 1-2002, item XBD/TC1/D6/7 is applied, adding the *ctime_r()* and
 2416 *localtime_r()* functions to the list of functions that use the *TZ* environment variable.

2417 A.9 Regular Expressions

2418 Rather than repeating the description of REs for each utility supporting REs, the standard
 2419 developers preferred a common, comprehensive description of regular expressions in one place.
 2420 The most common behavior is described here, and exceptions or extensions to this are
 2421 documented for the respective utilities, as appropriate.

2422 The BRE corresponds to the *ed* or historical *grep* type, and the ERE corresponds to the historical
 2423 *egrep* type (now *grep -E*).

2424 The text is based on the *ed* description and substantially modified, primarily to aid developers
 2425 and others in the understanding of the capabilities and limitations of REs. Much of this was
 2426 influenced by internationalization requirements.

2427 It should be noted that the definitions in this section do not cover the *tr* utility; the *tr* syntax does
 2428 not employ REs.

2429 The specification of REs is particularly important to internationalization because pattern
 2430 matching operations are very basic operations in business and other operations. The syntax and
 2431 rules of REs are intended to be as intuitive as possible to make them easy to understand and use.
 2432 The historical rules and behavior do not provide that capability to non-English language users,
 2433 and do not provide the necessary support for commonly used characters and language
 2434 constructs. It was necessary to provide extensions to the historical RE syntax and rules to
 2435 accommodate other languages.

2436 As they are limited to bracket expressions, the rationale for these modifications is in the Base
 2437 Definitions volume of IEEE Std 1003.1-200x, Section 9.3.5, RE Bracket Expression.

2438 A.9.1 Regular Expression Definitions

2439 It is possible to determine what strings correspond to subexpressions by recursively applying
 2440 the leftmost longest rule to each subexpression, but only with the proviso that the overall match
 2441 is leftmost longest. For example, matching "`\(ac*\)c*d[ac]*\1`" against `acdacaaa` matches
 2442 `acdacaaa` (with `\1=a`); simply matching the longest match for "`\(ac*\)`" would yield `\1=ac`, but
 2443 the overall match would be smaller (`acdac`). Conceptually, the implementation must examine
 2444 every possible match and among those that yield the leftmost longest total matches, pick the one
 2445 that does the longest match for the leftmost subexpression, and so on. Note that this means that
 2446 matching by subexpressions is context-dependent: a subexpression within a larger RE may
 2447 match a different string from the one it would match as an independent RE, and two instances of
 2448 the same subexpression within the same larger RE may match different lengths even in similar
 2449 sequences of characters. For example, in the ERE "`(a.*b)(a.*b)`", the two identical
 2450 subexpressions would match four and six characters, respectively, of `acbbaccccb`.

2451 The definition of single character has been expanded to include also collating elements
 2452 consisting of two or more characters; this expansion is applicable only when a bracket
 2453 expression is included in the BRE or ERE. An example of such a collating element may be the
 2454 Dutch *ij*, which collates as a 'y'. In some encodings, a ligature "i with j" exists as a character
 2455 and would represent a single-character collating element. In another encoding, no such ligature
 2456 exists, and the two-character sequence *ij* is defined as a multi-character collating element.
 2457 Outside brackets, the *ij* is treated as a two-character RE and matches the same characters in a
 2458 string. Historically, a bracket expression only matched a single character. The ISO POSIX-2:1993
 2459 standard required bracket expressions like "`[[:lower:]]`" to match multi-character collating
 2460 elements such as "*ij*". However, this requirement led to behavior that many users did not
 2461 expect and that could not feasibly be mimicked in user code, and it was rarely if ever
 2462 implemented correctly. The current standard leaves it unspecified whether a bracket expression
 2463 matches a multi-character collating element, allowing both historical and ISO POSIX-2:1993
 2464 standard implementations to conform.

2465 Also, in the current standard, it is unspecified whether character class expressions like
 2466 "`[[:lower:]]`" can include multi-character collating elements like "*ij*"; hence
 2467 "`[[:lower:]]`" can match "*ij*", and "`[^[:lower:]]`" can fail to match "*ij*". Common
 2468 practice is for a character class expression to match a collating element if it matches the collating
 2469 element's first character.

2470 A.9.2 Regular Expression General Requirements

2471 The definition of which sequence is matched when several are possible is based on the leftmost-
 2472 longest rule historically used by deterministic recognizers. This rule is easier to define and
 2473 describe, and arguably more useful, than the first-match rule historically used by non-
 2474 deterministic recognizers. It is thought that dependencies on the choice of rule are rare; carefully
 2475 contrived examples are needed to demonstrate the difference.

2476 A formal expression of the leftmost-longest rule is:

2477 The search is performed as if all possible suffixes of the string were tested for a prefix
 2478 matching the pattern; the longest suffix containing a matching prefix is chosen, and the
 2479 longest possible matching prefix of the chosen suffix is identified as the matching
 2480 sequence.

2481 Historically, most RE implementations only match lines, not strings. However, that is more an
 2482 effect of the usage than of an inherent feature of REs themselves. Consequently,
 2483 IEEE Std 1003.1-200x does not regard <newline>s as special; they are ordinary characters, and
 2484 both a period and a non-matching list can match them. Those utilities (like *grep*) that do not
 2485 allow <newline>s to match are responsible for eliminating any <newline> from strings before
 2486 matching against the RE. The *regcomp()* function, however, can provide support for such
 2487 processing without violating the rules of this section.

2488 Some implementations of *egrep* have had very limited flexibility in handling complex EREs.
 2489 IEEE Std 1003.1-200x does not attempt to define the complexity of a BRE or ERE, but does place
 2490 a lower limit on it—any RE must be handled, as long as it can be expressed in 256 bytes or less.
 2491 (Of course, this does not place an upper limit on the implementation.) There are historical
 2492 programs using a non-deterministic-recognizer implementation that should have no difficulty
 2493 with this limit. It is possible that a good approach would be to attempt to use the faster, but
 2494 more limited, deterministic recognizer for simple expressions and to fall back on the non-
 2495 deterministic recognizer for those expressions requiring it. Non-deterministic implementations
 2496 must be careful to observe the rules on which match is chosen; the longest match, not the first
 2497 match, starting at a given character is used.

2498 The term “invalid” highlights a difference between this section and some others:
 2499 IEEE Std 1003.1-200x frequently avoids mandating of errors for syntax violations because they
 2500 can be used by implementors to trigger extensions. However, the authors of the
 2501 internationalization features of REs wanted to mandate errors for certain conditions to identify
 2502 usage problems or non-portable constructs. These are identified within this rationale as
 2503 appropriate. The remaining syntax violations have been left implicitly or explicitly undefined.
 2504 For example, the BRE construct “\{1,2,3\}” does not comply with the grammar. A
 2505 conforming application cannot rely on it producing an error nor matching the literal characters
 2506 “\{1,2,3\}”.

2507 The term “undefined” was used in favor of “unspecified” because many of the situations are
 2508 considered errors on some implementations, and the standard developers considered that
 2509 consistency throughout the section was preferable to mixing undefined and unspecified.

2510 A.9.3 Basic Regular Expressions

2511 There is no additional rationale provided for this section.

2512 A.9.3.1 BREs Matching a Single Character or Collating Element

2513 There is no additional rationale provided for this section.

2514 A.9.3.2 BRE Ordinary Characters

2515 There is no additional rationale provided for this section.

2516 A.9.3.3 BRE Special Characters

2517 There is no additional rationale provided for this section.

2518 A.9.3.4 Periods in BREs

2519 There is no additional rationale provided for this section.

2520 A.9.3.5 RE Bracket Expression

2521 Range expressions are, historically, an integral part of REs. However, the requirements of
2522 “natural language behavior” and portability do conflict. In the POSIX locale, ranges must be
2523 treated according to the collating sequence and include such characters that fall within the range
2524 based on that collating sequence, regardless of character values. In other locales, ranges have
2525 unspecified behavior.

2526 Some historical implementations allow range expressions where the ending range point of one
2527 range is also the starting point of the next (for instance, "[a-m-o]"). This behavior should not
2528 be permitted, but to avoid breaking historical implementations, it is now *undefined* whether it is
2529 a valid expression and how it should be interpreted.

2530 Current practice in *awk* and *lex* is to accept escape sequences in bracket expressions as per the
2531 Base Definitions volume of IEEE Std 1003.1-200x, Table 5-1, Escape Sequences and Associated
2532 Actions, while the normal ERE behavior is to regard such a sequence as consisting of two
2533 characters. Allowing the *awk/lex* behavior in EREs would change the normal behavior in an
2534 unacceptable way; it is expected that *awk* and *lex* will decode escape sequences in EREs before
2535 passing them to *regcomp()* or comparable routines. Each utility describes the escape sequences it
2536 accepts as an exception to the rules in this section; the list is not the same, for historical reasons.

2537 As noted previously, the new syntax and rules have been added to accommodate other
2538 languages than English. The remainder of this section describes the rationale for these
2539 modifications.

2540 In the POSIX locale, a regular expression that starts with a range expression matches a set of
2541 strings that are contiguously sorted, but this is not necessarily true in other locales. For example,
2542 a French locale might have the following behavior:

```
2543 $ ls
2544 alpha Alpha estimé ESTIMÉ été eurêka
2545 $ ls [a-e]*
2546 alpha Alpha estimé eurêka
```

2547 Such disagreements between matching and contiguous sorting are unavoidable because POSIX
2548 sorting cannot be implemented in terms of a deterministic finite-state automaton (DFA), but
2549 range expressions by design are implementable in terms of DFAs.

2550 Historical implementations used native character order to interpret range expressions. The
2551 ISO POSIX-2:1993 standard instead required collating element order (CEO): the order that
2552 collating elements were specified between the **order_start** and **order_end** keywords in the
2553 **LC_COLLATE** category of the current locale. CEO had some advantages in portability over the

2554 native character order, but it also had some disadvantages:

- 2555 • CEO could not feasibly be mimicked in user code, leading to inconsistencies between
2556 POSIX matchers and matchers in popular user programs like Emacs, *ksh*, and Perl.
- 2557 • CEO caused range expressions to match accented and capitalized letters contrary to many
2558 users' expectations. For example, "[a-e]" typically matched both 'E' and 'á' but
2559 neither 'A' nor 'é'.
- 2560 • CEO was not consistent across implementations. In practice, CEO was often less portable
2561 than native character order. For example, it was common for the CEOs of two
2562 implementation-supplied locales to disagree, even if both locales were named "da_DK".

2563 Because of these problems, some implementations of regular expressions continued to use native
2564 character order. Others used the collation sequence, which is more consistent with sorting than
2565 either CEO or native order, but which departs further from the traditional POSIX semantics
2566 because it generally requires "[a-e]" to match either 'A' or 'E' but not both. As a result of
2567 this kind of implementation variation, programmers who wanted to write portable regular
2568 expressions could not rely on the ISO POSIX-2:1993 standard guarantees in practice.

2569 While revising the standard, lengthy consideration was given to proposals to attack this problem
2570 by adding an API for querying the CEO to allow user-mode matchers, but none of these
2571 proposals had implementation experience and none achieved consensus. Leaving the standard
2572 alone was also considered, but rejected due to the problems described above.

2573 The current standard leaves unspecified the behavior of a range expression outside the POSIX
2574 locale. This makes it clearer that conforming applications should avoid range expressions
2575 outside the POSIX locale, and it allows implementations and compatible user-mode matchers to
2576 interpret range expressions using native order, CEO, collation sequence, or other, more
2577 advanced techniques. The concerns which led to this change were raised in IEEE PASC
2578 interpretation 1003.2 #43 and others, and related to ambiguities in the specification of how
2579 multi-character collating elements should be handled in range expressions. These ambiguities
2580 had led to multiple interpretations of the specification, in conflicting ways, which led to varying
2581 implementations. As noted above, efforts were made to resolve the differences, but no solution
2582 has been found that would be specific enough to allow for portable software while not
2583 invalidating existing implementations.

2584 The standard developers recognize that collating elements are important, such elements being
2585 common in several European languages; for example, 'ch' or 'll' in traditional Spanish;
2586 'aa' in several Scandinavian languages. Existing internationalized implementations have
2587 processed, and continue to process, these elements in range expressions. Efforts are expected to
2588 continue in the future to find a way to define the behavior of these elements precisely and
2589 portably.

2590 The ISO POSIX-2:1993 standard required "[b-a]" to be an invalid expression in the POSIX
2591 locale, but this requirement has been relaxed in this version of the standard so that "[b-a]" can
2592 instead be treated as a valid expression that does not match any string.

2593 A.9.3.6 BREs Matching Multiple Characters

2594 The limit of nine back-references to subexpressions in the RE is based on the use of a single-digit
2595 identifier; increasing this to multiple digits would break historical applications. This does not
2596 imply that only nine subexpressions are allowed in REs. The following is a valid BRE with ten
2597 subexpressions:

```
2598 \(\(\(\(ab\) *c\) *d\) \(\(ef\) * \(\(gh\) \{2\} \(\(ij\) * \(\(kl\) * \(\(mn\) * \(\(op\) * \(\(qr\) *
```

2599 The standard developers regarded the common historical behavior, which supported "\n*", but
2600 not "\n\{min,max\}", "\(\dots\) *", or "\(\dots\) \{min,max\}", as a non-intentional

2601 result of a specific implementation, and they supported both duplication and interval
2602 expressions following subexpressions and back-references.

2603 The changes to the processing of the back-reference expression remove an unspecified or
2604 ambiguous behavior in the Shell and Utilities volume of IEEE Std 1003.1-200x, aligning it with
2605 the requirements specified for the *regcomp()* expression, and is the result of PASC Interpretation
2606 1003.2-92 #43 submitted for the ISO POSIX-2: 1993 standard.

2607 A.9.3.7 BRE Precedence

2608 There is no additional rationale provided for this section.

2609 A.9.3.8 BRE Expression Anchoring

2610 Often, the dollar sign is viewed as matching the ending <newline> in text files. This is not
2611 strictly true; the <newline> is typically eliminated from the strings to be matched, and the dollar
2612 sign matches the terminating null character.

2613 The ability of '^', '\$', and '*' to be non-special in certain circumstances may be confusing to
2614 some programmers, but this situation was changed only in a minor way from historical practice
2615 to avoid breaking many historical scripts. Some consideration was given to making the use of
2616 the anchoring characters undefined if not escaped and not at the beginning or end of strings.
2617 This would cause a number of historical BREs, such as "2^10", "\$HOME", and "\$1.35", that
2618 relied on the characters being treated literally, to become invalid.

2619 However, one relatively uncommon case was changed to allow an extension used on some
2620 implementations. Historically, the BREs "f^oo" and "\(^foo\)" did not match the same
2621 string, despite the general rule that subexpressions and entire BREs match the same strings. To
2622 increase consensus, IEEE Std 1003.1-200x has allowed an extension on some implementations to
2623 treat these two cases in the same way by declaring that anchoring *may* occur at the beginning or
2624 end of a subexpression. Therefore, portable BREs that require a literal circumflex at the
2625 beginning or a dollar sign at the end of a subexpression must escape them. Note that a BRE such
2626 as "a\(^bc\)" will either match "a^bc" or nothing on different systems under the rules.

2627 ERE anchoring has been different from BRE anchoring in all historical systems. An unescaped
2628 anchor character has never matched its literal counterpart outside a bracket expression. Some
2629 implementations treated "foo\$bar" as a valid expression that never matched anything; others
2630 treated it as invalid. IEEE Std 1003.1-200x mandates the former, valid unmatched behavior.

2631 Some implementations have extended the BRE syntax to add alternation. For example, the
2632 subexpression "\ (foo\$ | bar\)" would match either "foo" at the end of the string or "bar"
2633 anywhere. The extension is triggered by the use of the undefined "\|" sequence. Because the
2634 BRE is undefined for portable scripts, the extending system is free to make other assumptions,
2635 such that the '\$' represents the end-of-line anchor in the middle of a subexpression. If it were
2636 not for the extension, the '\$' would match a literal dollar sign under the rules.

2637 A.9.4 Extended Regular Expressions

2638 As with BREs, the standard developers decided to make the interpretation of escaped ordinary
2639 characters undefined.

2640 The right parenthesis is not listed as an ERE special character because it is only special in the
2641 context of a preceding left parenthesis. If found without a preceding left parenthesis, the right
2642 parenthesis has no special meaning.

2643 The interval expression, "{m,n}", has been added to EREs. Historically, the interval expression
2644 has only been supported in some ERE implementations. The standard developers estimated that
2645 the addition of interval expressions to EREs would not decrease consensus and would also make
2646 BREs more of a subset of EREs than in many historical implementations.

2647 It was suggested that, in addition to interval expressions, back-references ('\`\n`') should also be
2648 added to EREs. This was rejected by the standard developers as likely to decrease consensus.

2649 In historical implementations, multiple duplication symbols are usually interpreted from left to
2650 right and treated as additive. As an example, "`a+*b`" matches zero or more instances of '`a`'
2651 followed by a '`b`'. In IEEE Std 1003.1-200x, multiple duplication symbols are undefined; that is,
2652 they cannot be relied upon for conforming applications. One reason for this is to provide some
2653 scope for future enhancements.

2654 The precedence of operations differs between EREs and those in *lex*; in *lex*, for historical reasons,
2655 interval expressions have a lower precedence than concatenation.

2656 A.9.4.1 EREs Matching a Single Character or Collating Element

2657 There is no additional rationale provided for this section.

2658 A.9.4.2 ERE Ordinary Characters

2659 There is no additional rationale provided for this section.

2660 A.9.4.3 ERE Special Characters

2661 There is no additional rationale provided for this section.

2662 A.9.4.4 Periods in EREs

2663 There is no additional rationale provided for this section.

2664 A.9.4.5 ERE Bracket Expression

2665 There is no additional rationale provided for this section.

2666 A.9.4.6 EREs Matching Multiple Characters

2667 There is no additional rationale provided for this section.

2668 A.9.4.7 ERE Alternation

2669 There is no additional rationale provided for this section.

2670 A.9.4.8 ERE Precedence

2671 There is no additional rationale provided for this section.

2672 A.9.4.9 ERE Expression Anchoring

2673 There is no additional rationale provided for this section.

2674 A.9.5 Regular Expression Grammar

2675 The grammars are intended to represent the range of acceptable syntaxes available to
2676 conforming applications. There are instances in the text where undefined constructs are
2677 described; as explained previously, these allow implementation extensions. There is no intended
2678 requirement that an implementation extension must somehow fit into the grammars shown
2679 here.

2680 The BRE grammar does not permit `L_ANCHOR` or `R_ANCHOR` inside "`\(`" and "`\)`" (which
2681 implies that '`^`' and '`$`' are ordinary characters). This reflects the semantic limits on the
2682 application, as noted in the Base Definitions volume of IEEE Std 1003.1-200x, Section 9.3.8, BRE
2683 Expression Anchoring. Implementations are permitted to extend the language to interpret '`^`'
2684 and '`$`' as anchors in these locations, and as such, conforming applications cannot use
2685 unescaped '`^`' and '`$`' in positions inside "`\(`" and "`\)`" that might be interpreted as
2686 anchors.

2687 The ERE grammar does not permit several constructs that the Base Definitions volume of

2688 IEEE Std 1003.1-200x, Section 9.4.2, ERE Ordinary Characters and the Base Definitions volume of
2689 IEEE Std 1003.1-200x, Section 9.4.3, ERE Special Characters specify as having undefined results:

- 2690 • ORD_CHAR preceded by '\'
- 2691 • ERE_dupl_symbol(s) appearing first in an ERE, or immediately following '|', '^', or '('
- 2692 • '{' not part of a valid ERE_dupl_symbol
- 2693 • '|' appearing first or last in an ERE, or immediately following '|' or '(', or
2694 immediately preceding ')'

2695 Implementations are permitted to extend the language to allow these. Conforming applications
2696 cannot use such constructs.

2697 A.9.5.1 BRE/ERE Grammar Lexical Conventions

2698 There is no additional rationale provided for this section.

2699 A.9.5.2 RE and Bracket Expression Grammar

2700 The removal of the *Back_open_paren Back_close_paren* option from the *nondupl_RE* specification is
2701 the result of PASC Interpretation 1003.2-92 #43 submitted for the ISO POSIX-2:1993 standard.
2702 Although the grammar required support for null subexpressions, this section does not describe
2703 the meaning of, and historical practice did not support, this construct.

2704 A.9.5.3 ERE Grammar

2705 There is no additional rationale provided for this section.

2706 A.10 Directory Structure and Devices

2707 A.10.1 Directory Structure and Files

2708 A description of the historical **/usr/tmp** was omitted, removing any concept of differences in
2709 emphasis between the **/** and **/usr** directories. The descriptions of **/bin**, **/usr/bin**, **/lib**, and **/usr/lib**
2710 were omitted because they are not useful for applications. In an early draft, a distinction was
2711 made between system and application directory usage, but this was not found to be useful.

2712 The directories **/** and **/dev** are included because the notion of a hierarchical directory structure is
2713 key to other information presented elsewhere in IEEE Std 1003.1-200x. In early drafts, it was
2714 argued that special devices and temporary files could conceivably be handled without a
2715 directory structure on some implementations. For example, the system could treat the characters
2716 **"/tmp"** as a special token that would store files using some non-POSIX file system structure.
2717 This notion was rejected by the standard developers, who required that all the files in this
2718 section be implemented via POSIX file systems.

2719 The **/tmp** directory is retained in IEEE Std 1003.1-200x to accommodate historical applications
2720 that assume its availability. Implementations are encouraged to provide suitable directory names
2721 in the environment variable **TMPDIR** and applications are encouraged to use the contents of
2722 **TMPDIR** for creating temporary files.

2723 The standard files **/dev/null** and **/dev/tty** are required to be both readable and writable to allow
2724 applications to have the intended historical access to these files.

2725 The standard file **/dev/console** has been added for alignment with the Single UNIX
2726 Specification.

2727
2728
2729

A.10.2 Output Devices and Terminal Types

IEEE Std 1003.1-2001/Cor 2-2004, item XBD/TC2/D6/17 is applied, making it clear that the requirements for documenting terminal support are in the system documentation.

2730
2731
2732
2733
2734
2735
2736
2737
2738
2739

A.11 General Terminal Interface

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

2740
2741
2742
2743
2744
2745
2746
2747
2748

Although the needs of both interface implementors and application developers were addressed throughout IEEE Std 1003.1-200x, this section pays more attention to the needs of the latter. This is because, while many aspects of the programming interface can be hidden from the user by the application developer, the terminal interface is usually a large part of the user interface. Although to some extent the application developer can build missing features or work around inappropriate ones, the difficulties of doing that are greater in the terminal interface than elsewhere. For example, efficiency prohibits the average program from interpreting every character passing through it in order to simulate character erase, line kill, and so on. These functions should usually be done by the operating system, possibly at the interrupt level.

2749
2750
2751
2752
2753
2754

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

2755
2756
2757
2758
2759
2760

The advantages of this new method include:

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

2761
2762
2763
2764
2765
2766

The disadvantages include:

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

The issue of modem control was excluded from IEEE Std 1003.1-200x on the grounds that:

- It was concerned with setting and control of hardware timers.

- 2767 • The appropriate timers and settings vary widely internationally.
- 2768 • Feedback from European computer manufacturers indicated that this facility was not
- 2769 consistent with European needs and that specification of such a facility was not a
- 2770 requirement for portability.

2771 A.11.1 Interface Characteristics

2772 A.11.1.1 Opening a Terminal Device File

2773 There is no additional rationale provided for this section.

2774 A.11.1.2 Process Groups

2775 There is a potential race when the members of the foreground process group on a terminal leave
 2776 that process group, either by exit or by changing process groups. After the last process exits the
 2777 process group, but before the foreground process group ID of the terminal is changed (usually
 2778 by a job control shell), it would be possible for a new process to be created with its process ID
 2779 equal to the terminal's foreground process group ID. That process might then become the
 2780 process group leader and accidentally be placed into the foreground on a terminal that was not
 2781 necessarily its controlling terminal. As a result of this problem, the controlling terminal is
 2782 defined to not have a foreground process group during this time.

2783 The cases where a controlling terminal has no foreground process group occur when all
 2784 processes in the foreground process group either terminate and are waited for or join other
 2785 process groups via *setpgid()* or *setsid()*. If the process group leader terminates, this is the first
 2786 case described; if it leaves the process group via *setpgid()*, this is the second case described (a
 2787 process group leader cannot successfully call *setsid()*). When one of those cases causes a
 2788 controlling terminal to have no foreground process group, it has two visible effects on
 2789 applications. The first is the value returned by *tcgetpgrp()*. The second (which occurs only in the
 2790 case where the process group leader terminates) is the sending of signals in response to special
 2791 input characters. The intent of IEEE Std 1003.1-200x is that no process group be wrongly
 2792 identified as the foreground process group by *tcgetpgrp()* or unintentionally receive signals
 2793 because of placement into the foreground.

2794 In 4.3 BSD, the old process group ID continues to be used to identify the foreground process
 2795 group and is returned by the function equivalent to *tcgetpgrp()*. In that implementation it is
 2796 possible for a newly created process to be assigned the same value as a process ID and then form
 2797 a new process group with the same value as a process group ID. The result is that the new
 2798 process group would receive signals from this terminal for no apparent reason, and IEEE Std
 2799 1003.1-200x precludes this by forbidding a process group from entering the foreground
 2800 in this way. It would be more direct to place part of the requirement made by the last sentence
 2801 under *fork()*, but there is no convenient way for that section to refer to the value that *tcgetpgrp()*
 2802 returns, since in this case there is no process group and thus no process group ID.

2803 One possibility for a conforming implementation is to behave similarly to 4.3 BSD, but to
 2804 prevent this reuse of the ID, probably in the implementation of *fork()*, as long as it is in use by
 2805 the terminal.

2806 Another possibility is to recognize when the last process stops using the terminal's foreground
 2807 process group ID, which is when the process group lifetime ends, and to change the terminal's
 2808 foreground process group ID to a reserved value that is never used as a process ID or process
 2809 group ID. (See the definition of *process group lifetime* in the definitions section.) The process ID
 2810 can then be reserved until the terminal has another foreground process group.

2811 The 4.3 BSD implementation permits the leader (and only member) of the foreground process
 2812 group to leave the process group by calling the equivalent of *setpgid()* and to later return,

2813 expecting to return to the foreground. There are no known application needs for this behavior,
 2814 and IEEE Std 1003.1-200x neither requires nor forbids it (except that it is forbidden for session
 2815 leaders) by leaving it unspecified.

2816 A.11.1.3 *The Controlling Terminal*

2817 IEEE Std 1003.1-200x does not specify a mechanism by which to allocate a controlling terminal.
 2818 This is normally done by a system utility (such as *getty*) and is considered an administrative
 2819 feature outside the scope of IEEE Std 1003.1-200x.

2820 Historical implementations allocate controlling terminals on certain *open()* calls. Since *open()* is
 2821 part of POSIX.1, its behavior had to be dealt with. The traditional behavior is not required
 2822 because it is not very straightforward or flexible for either implementations or applications.
 2823 However, because of its prevalence, it was not practical to disallow this behavior either. Thus, a
 2824 mechanism was standardized to ensure portable, predictable behavior in *open()*.

2825 Some historical implementations deallocate a controlling terminal on the last system-wide close.
 2826 This behavior is neither required nor prohibited. Even on implementations that do provide this
 2827 behavior, applications generally cannot depend on it due to its system-wide nature.

2828 A.11.1.4 *Terminal Access Control*

2829 The access controls described in this section apply only to a process that is accessing its
 2830 controlling terminal. A process accessing a terminal that is not its controlling terminal is
 2831 effectively treated the same as a member of the foreground process group. While this may seem
 2832 unintuitive, note that these controls are for the purpose of job control, not security, and job
 2833 control relates only to the controlling terminal of a process. Normal file access permissions
 2834 handle security.

2835 If the process calling *read()* or *write()* is in a background process group that is orphaned, it is not
 2836 desirable to stop the process group, as it is no longer under the control of a job control shell that
 2837 could put it into the foreground again. Accordingly, calls to *read()* or *write()* functions by such
 2838 processes receive an immediate error return. This is different from 4.2 BSD, which kills orphaned
 2839 processes that receive terminal stop signals.

2840 The foreground/background/orphaned process group check performed by the terminal driver
 2841 must be repeatedly performed until the calling process moves into the foreground or until the
 2842 process group of the calling process becomes orphaned. That is, when the terminal driver
 2843 determines that the calling process is in the background and should receive a job control signal,
 2844 it sends the appropriate signal (SIGTTIN or SIGTTOU) to every process in the process group of
 2845 the calling process and then it allows the calling process to immediately receive the signal. The
 2846 latter is typically performed by blocking the process so that the signal is immediately noticed.
 2847 Note, however, that after the process finishes receiving the signal and control is returned to the
 2848 driver, the terminal driver must re-execute the foreground/background/orphaned process
 2849 group check. The process may still be in the background, either because it was continued in the
 2850 background by a job control shell, or because it caught the signal and did nothing.

2851 The terminal driver repeatedly performs the foreground/background/orphaned process group
 2852 checks whenever a process is about to access the terminal. In the case of *write()* or the control
 2853 *tc*()* functions, the check is performed at the entry of the function. In the case of *read()*,
 2854 the check is performed not only at the entry of the function, but also after blocking the process to
 2855 wait for input characters (if necessary). That is, once the driver has determined that the process
 2856 calling the *read()* function is in the foreground, it attempts to retrieve characters from the input
 2857 queue. If the queue is empty, it blocks the process waiting for characters. When characters are
 2858 available and control is returned to the driver, the terminal driver must return to the repeated
 2859 foreground/background/orphaned process group check again. The process may have moved
 2860 from the foreground to the background while it was blocked waiting for input characters.

2861 A.11.1.5 *Input Processing and Reading Data*

2862 There is no additional rationale provided for this section.

2863 A.11.1.6 *Canonical Mode Input Processing*2864 The term “character” is intended here. ERASE should erase the last character, not the last byte.
2865 In the case of multi-byte characters, these two may be different.2866 4.3 BSD has a WERASE character that erases the last “word” typed (but not any preceding
2867 <blank>s or <tab>s). A word is defined as a sequence of non-<blank>s, with <tab>s counted as
2868 <blank>s. Like ERASE, WERASE does not erase beyond the beginning of the line. This
2869 WERASE feature has not been specified in POSIX.1 because it is difficult to define in the
2870 international environment. It is only useful for languages where words are delimited by
2871 <blank>s. In some ideographic languages, such as Japanese and Chinese, words are not
2872 delimited at all. The WERASE character should presumably go back to the beginning of a
2873 sentence in those cases; practically, this means it would not be used much for those languages.2874 It should be noted that there is a possible inherent deadlock if the application and
2875 implementation conflict on the value of {MAX_CANON}. With ICANON set (if IXOFF is
2876 enabled) and more than {MAX_CANON} characters transmitted without a <linefeed>,
2877 transmission will be stopped, the <linefeed> (or <carriage-return> when ICRLF is set) will never
2878 arrive, and the *read()* will never be satisfied.2879 An application should not set IXOFF if it is using canonical mode unless it knows that (even in
2880 the face of a transmission error) the conditions described previously cannot be met or unless it
2881 is prepared to deal with the possible deadlock in some other way, such as timeouts.2882 It should also be noted that this can be made to happen in non-canonical mode if the trigger
2883 value for sending IXOFF is less than VMIN and VTIME is zero.2884 A.11.1.7 *Non-Canonical Mode Input Processing*

2885 Some points to note about MIN and TIME:

- 2886 1. The interactions of MIN and TIME are not symmetric. For example, when MIN>0 and
-
- 2887 TIME=0, TIME has no effect. However, in the opposite case where MIN=0 and TIME>0,
-
- 2888 both MIN and TIME play a role in that MIN is satisfied with the receipt of a single
-
- 2889 character.
-
- 2890 2. Also note that in case A (MIN>0, TIME>0), TIME represents an inter-character timer,
-
- 2891 while in case C (MIN=0, TIME>0), TIME represents a read timer.

2892 These two points highlight the dual purpose of the MIN/TIME feature. Cases A and B, where
2893 MIN>0, exist to handle burst-mode activity (for example, file transfer programs) where a
2894 program would like to process at least MIN characters at a time. In case A, the inter-character
2895 timer is activated by a user as a safety measure; in case B, it is turned off.2896 Cases C and D exist to handle single-character timed transfers. These cases are readily adaptable
2897 to screen-based applications that need to know if a character is present in the input queue before
2898 refreshing the screen. In case C, the read is timed; in case D, it is not.2899 Another important note is that MIN is always just a minimum. It does not denote a record
2900 length. That is, if a program does a read of 20 bytes, MIN is 10, and 25 characters are present, 20
2901 characters are returned to the user. In the special case of MIN=0, this still applies: if more than
2902 one character is available, they all will be returned immediately.

2903 *A.11.1.8 Writing Data and Output Processing*

2904 There is no additional rationale provided for this section.

2905 *A.11.1.9 Special Characters*

2906 There is no additional rationale provided for this section.

2907 *A.11.1.10 Modem Disconnect*

2908 There is no additional rationale provided for this section.

2909 *A.11.1.11 Closing a Terminal Device File*2910 IEEE Std 1003.1-200x does not specify that a *close()* on a terminal device file include the
2911 equivalent of a call to *tcflow(fd,TCOON)*.2912 An implementation that discards output at the time *close()* is called after reporting the return
2913 value to the *write()* call that data was written does not conform with IEEE Std 1003.1-200x. An
2914 application has functions such as *tcdrain()*, *tcflush()*, and *tcflow()* available to obtain the detailed
2915 behavior it requires with respect to flushing of output.2916 At the time of the last close on a terminal device, an application relinquishes any ability to exert
2917 flow control via *tcflow()*.2918 **A.11.2 Parameters that Can be Set**2919 *A.11.2.1 The termios Structure*2920 This structure is part of an interface that, in general, retains the historic grouping of flags.
2921 Although a more optimal structure for implementations may be possible, the degree of change
2922 to applications would be significantly larger.2923 *A.11.2.2 Input Modes*2924 Some historical implementations treated a long break as multiple events, as many as one per
2925 character time. The wording in POSIX.1 explicitly prohibits this.2926 Although the ISTRIP flag is normally superfluous with today's terminal hardware and software,
2927 it is historically supported. Therefore, applications may be using ISTRIP, and there is no
2928 technical problem with supporting this flag. Also, applications may wish to receive only 7-bit
2929 input bytes and may not be connected directly to the hardware terminal device (for example,
2930 when a connection traverses a network).2931 Also, there is no requirement in general that the terminal device ensures that high-order bits
2932 beyond the specified character size are cleared. ISTRIP provides this function for 7-bit
2933 characters, which are common.2934 In dealing with multi-byte characters, the consequences of a parity error in such a character, or
2935 in an escape sequence affecting the current character set, are beyond the scope of POSIX.1 and
2936 are best dealt with by the application processing the multi-byte characters.2937 *A.11.2.3 Output Modes*2938 POSIX.1 does not describe postprocessing of output to a terminal or detailed control of that from
2939 a conforming application. (That is, translation of <newline> to <carriage-return> followed by
2940 <linefeed> or <tab> processing.) There is nothing that a conforming application should do to its
2941 output for a terminal because that would require knowledge of the operation of the terminal. It
2942 is the responsibility of the operating system to provide postprocessing appropriate to the output
2943 device, whether it is a terminal or some other type of device.

2944 Extensions to POSIX.1 to control the type of postprocessing already exist and are expected to
 2945 continue into the future. The control of these features is primarily to adjust the interface between
 2946 the system and the terminal device so the output appears on the display correctly. This should
 2947 be set up before use by any application.

2948 In general, both the input and output modes should not be set absolutely, but rather modified
 2949 from the inherited state.

2950 A.11.2.4 Control Modes

2951 This section could be misread that the symbol "CSIZE" is a title in the **termios** *c_flag* field.
 2952 Although it does serve that function, it is also a required symbol, as a literal reading of POSIX.1
 2953 (and the caveats about typography) would indicate.

2954 A.11.2.5 Local Modes

2955 Non-canonical mode is provided to allow fast bursts of input to be read efficiently while still
 2956 allowing single-character input.

2957 The ECHONL function historically has been in many implementations. Since there seems to be
 2958 no technical problem with supporting ECHONL, it is included in POSIX.1 to increase consensus.

2959 The alternate behavior possible when ECHOK or ECHOE are specified with ICANON is
 2960 permitted as a compromise depending on what the actual terminal hardware can do. Erasing
 2961 characters and lines is preferred, but is not always possible.

2962 A.11.2.6 Special Control Characters

2963 Permitting VMIN and VTIME to overlap with VEOF and VEOL was a compromise for historical
 2964 implementations. Only when backwards-compatibility of object code is a serious concern to an
 2965 implementor should an implementation continue this practice. Correct applications that work
 2966 with the overlap (at the source level) should also work if it is not present, but not the reverse.

2967 A.12 Utility Conventions

2968 A.12.1 Utility Argument Syntax

2969 The standard developers considered that recent trends toward diluting the SYNOPSIS sections
 2970 of historical reference pages to the equivalent of:

2971 `command [options][operands]`

2972 were a disservice to the reader. Therefore, considerable effort was placed into rigorous
 2973 definitions of all the command line arguments and their interrelationships. The relationships
 2974 depicted in the synopses are normative parts of IEEE Std 1003.1-200x; this information is
 2975 sometimes repeated in textual form, but that is only for clarity within context.

2976 The use of "undefined" for conflicting argument usage and for repeated usage of the same
 2977 option is meant to prevent conforming applications from using conflicting arguments or
 2978 repeated options unless specifically allowed (as is the case with *ls*, which allows simultaneous,
 2979 repeated use of the *-C*, *-l*, and *-1* options). Many historical implementations will tolerate this
 2980 usage, choosing either the first or the last applicable argument. This tolerance can continue, but
 2981 conforming applications cannot rely upon it. (Other implementations may choose to print usage
 2982 messages instead.)

2983 The use of "undefined" for conflicting argument usage also allows an implementation to make
 2984 reasonable extensions to utilities where the implementor considers mutually-exclusive options

2985 according to IEEE Std 1003.1-200x to have a sensible meaning and result.

2986 IEEE Std 1003.1-200x does not define the result of a command when an option-argument or
 2987 operand is not followed by ellipses and the application specifies more than one of that option-
 2988 argument or operand. This allows an implementation to define valid (although non-standard)
 2989 behavior for the utility when more than one such option or operand is specified.

2990 The requirements for option-arguments are summarized as follows:

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

2996 Earlier versions of this standard included obsolescent syntax which showed some options with
 2997 (mandatory) adjacent option-arguments in the SYNOPSIS for some utilities. These have since
 2998 been removed. For all options with mandatory option-arguments, the SYNOPSIS now shows
 2999 <blank>s between the option and the option-argument; however, historical usage has not been
 3000 consistent in this area; therefore, <blank>s are required to be used by conforming applications
 3001 and to be handled by all implementations, but implementations are also required to handle an
 3002 adjacent option-argument in order to preserve backwards-compatibility for old scripts. One of
 3003 the justifications for selecting the multiple-argument method was that the single-argument case
 3004 is inherently ambiguous when the option-argument can legitimately be a null string.

3005 IEEE Std 1003.1-200x explicitly states that digits are permitted as operands and option-
 3006 arguments. The lower and upper bounds for the values of the numbers used for operands and
 3007 option-arguments were derived from the ISO C standard values for {LONG_MIN} and
 3008 {LONG_MAX}. The requirement on the standard utilities is that numbers in the specified range
 3009 do not cause a syntax error, although the specification of a number need not be semantically
 3010 correct for a particular operand or option-argument of a utility. For example, the specification of:

3011 `dd obs=3000000000`

3012 would yield undefined behavior for the application and could be a syntax error because the
 3013 number 3 000 000 000 is outside of the range -2 147 483 647 to +2 147 483 647. On the other hand:

3014 `dd obs=2000000000`

3015 may cause some error, such as “blocksize too large”, rather than a syntax error.

3016 A.12.2 Utility Syntax Guidelines

3017 This section is based on the rules listed in the SVID. It was included for two reasons:

- 3018 1. The individual utility descriptions in the Shell and Utilities volume of
 3019 IEEE Std 1003.1-200x, Chapter 4, Utilities needed a set of common (although not
 3020 universal) actions on which they could anchor their descriptions of option and operand
 3021 syntax. Most of the standard utilities actually do use these guidelines, and many of their
 3022 historical implementations use the *getopt()* function for their parsing. Therefore, it was
 3023 simpler to cite the rules and merely identify exceptions.
- 3024 2. Writers of conforming applications need suggested guidelines if the POSIX community is
 3025 to avoid the chaos of historical UNIX system command syntax.

3026 It is recommended that all *future* utilities and applications use these guidelines to enhance “user
 3027 portability”. The fact that some historical utilities could not be changed (to avoid breaking
 3028 historical applications) should not deter this future goal.

3029 The voluntary nature of the guidelines is highlighted by repeated uses of the word *should*
 3030 throughout. This usage should not be misinterpreted to imply that utilities that claim
 3031 conformance in their OPTIONS sections do not always conform.

3032 Guidelines 1 and 2 encourage utility writers to use only characters from the portable character
 3033 set because use of locale-specific characters may make the utility inaccessible from other locales.
 3034 Use of uppercase letters is discouraged due to problems associated with porting utilities to
 3035 systems that do not distinguish between uppercase and lowercase characters in filenames. Use
 3036 of non-alphanumeric characters is discouraged due to the number of utilities that treat non-
 3037 alphanumeric characters in “special” ways depending on context (such as the shell using
 3038 whitespace characters to delimit arguments, various quote characters for quoting, the dollar sign
 3039 to introduce variable expansion, etc.).

3040 In the Shell and Utilities volume of IEEE Std 1003.1-200x, Section 2.9.1, Simple Commands, it is
 3041 further stated that a command used in the Shell Command Language cannot be named with a
 3042 trailing colon.

3043 Guideline 3 was changed to allow alphanumeric characters (letters and digits) from the
 3044 character set to allow compatibility with historical usage. Historical practice allows the use of
 3045 digits wherever practical, and there are no portability issues that would prohibit the use of
 3046 digits. In fact, from an internationalization viewpoint, digits (being non-language-dependent)
 3047 are preferable over letters (a `-2` is intuitively self-explanatory to any user, while in the `-f filename`
 3048 the letter ‘f’ is a mnemonic aid only to speakers of Latin-based languages where “filename”
 3049 happens to translate to a word that begins with ‘f’). Since Guideline 3 still retains the word
 3050 “single”, multi-digit options are not allowed. Instances of historical utilities that used them have
 3051 been marked obsolescent, with the numbers being changed from option names to option-
 3052 arguments.

3053 It was difficult to achieve a satisfactory solution to the problem of name space in option
 3054 characters. When the standard developers desired to extend the historical `cc` utility to accept
 3055 ISO C standard programs, they found that all of the portable alphabet was already in use by
 3056 various vendors. Thus, they had to devise a new name, `c89` (now superseded by `c99`), rather than
 3057 something like `cc -X`. There were suggestions that implementors be restricted to providing
 3058 extensions through various means (such as using a plus sign as the option delimiter or using
 3059 option characters outside the alphanumeric set) that would reserve all of the remaining
 3060 alphanumeric characters for future POSIX standards. These approaches were resisted because
 3061 they lacked the historical style of UNIX systems. Furthermore, if a vendor-provided option
 3062 should become commonly used in the industry, it would be a candidate for standardization. It
 3063 would be desirable to standardize such a feature using historical practice for the syntax (the
 3064 semantics can be standardized with any syntax). This would not be possible if the syntax was
 3065 one reserved for the vendor. However, since the standardization process may lead to minor
 3066 changes in the semantics, it may prove to be better for a vendor to use a syntax that will not be
 3067 affected by standardization.

3068 Guideline 8 includes the concept of comma-separated lists in a single argument. It is up to the
 3069 utility to parse such a list itself because `getopt()` just returns the single string. This situation was
 3070 retained so that certain historical utilities would not violate the guidelines. Applications
 3071 preparing for international use should be aware of an occasional problem with comma-
 3072 separated lists: in some locales, the comma is used as the radix character. Thus, if an application
 3073 is preparing operands for a utility that expects a comma-separated list, it should avoid
 3074 generating non-integer values through one of the means that is influenced by setting the
 3075 `LC_NUMERIC` variable (such as `awk`, `bc`, `printf`, or `printf()`).

3076 Unless explicitly stated otherwise in the utility description, Guideline 9 requires applications to
 3077 put options before operands, and requires utilities to accept any such usage without
 3078 misinterpreting operands as options. For example, if an implementation of the `printf` utility

3079 supports a `-e` option as an extension, the command:

```
3080 printf %s -e
```

3081 must output the string `"-e"` without interpreting the `-e` as an option.

3082 Applications calling any utility with a first operand starting with `'-'` should usually specify `--`,
 3083 as indicated by Guideline 10, to mark the end of the options. This is true even if the SYNOPSIS
 3084 in the Shell and Utilities volume of IEEE Std 1003.1-200x does not specify any options;
 3085 implementations may provide options as extensions to the Shell and Utilities volume of
 3086 IEEE Std 1003.1-200x. The standard utilities that do not support Guideline 10 indicate that fact in
 3087 the OPTIONS section of the utility description.

3088 Guideline 7 allows any string to be an option-argument; an option-argument can begin with any
 3089 character, can be `-` or `--`, and can be an empty string. For example, the commands `pr -h -`, `pr -h`
 3090 `--`, `pr -h -d`, `pr -h +2`, and `pr -h ''` contain the option-arguments `-`, `--`, `-d`, `+2`, and an empty
 3091 string, respectively. Conversely, the command `pr -h -- -d` treats `-d` as an option, not as an
 3092 argument, because the `--` is an option-argument here, not a delimiter.

3093 Guideline 11 was modified to clarify that the order of different options should not matter
 3094 relative to one another. However, the order of repeated options that also have option-arguments
 3095 may be significant; therefore, such options are required to be interpreted in the order that they
 3096 are specified. The `make` utility is an instance of a historical utility that uses repeated options in
 3097 which the order is significant. Multiple files are specified by giving multiple instances of the `-f`
 3098 option; for example:

```
3099 make -f common_header -f specific_rules target
```

3100 Guideline 13 does not imply that all of the standard utilities automatically accept the operand
 3101 `'-'` to mean standard input or output, nor does it specify the actions of the utility upon
 3102 encountering multiple `'-'` operands. It simply says that, by default, `'-'` operands are not used
 3103 for other purposes in the file reading or writing (but not when using `stat()`, `unlink()`, `touch`, and
 3104 so on) utilities. In previous revisions of this standard, all information concerning actual
 3105 treatment of the `'-'` operand is found in the individual utility sections. Many implementations,
 3106 however, treated `'-'` as standard input or output and many applications depended on this
 3107 behavior even though it was not standard. This behavior is now implementation-defined.
 3108 Portable applications should not use `'-'` to mean standard input or output unless it is explicitly
 3109 stated to do so in the utility description and they should always use `./-` if they intend to refer
 3110 to a file named `'-'` in the current working directory.

3111 Guideline 14 is intended to prohibit implementations that would treat the command `ls -l -d` as if
 3112 it were `ls -- -l -d` or `ls -l -- -d`.

3113 The standard permits implementations to have extensions that violate the Utility Syntax
 3114 Guidelines so long as when the utility is used in line with the forms defined by the standard it
 3115 follows the Utility Syntax Guidelines. Thus, `head-42file` and `ls--help` are permitted extensions.
 3116 The intent is to allow extensions so long as the standard form is accepted and follows the
 3117 guidelines.

3118 An area of concern was that as implementations mature, implementation-defined utilities and
 3119 implementation-defined utility options will result. The idea was expressed that there needed to
 3120 be a standard way, say an environment variable or some such mechanism, to identify
 3121 implementation-defined utilities separately from standard utilities that may have the same
 3122 name. It was decided that there already exist several ways of dealing with this situation and that
 3123 it is outside of the scope to attempt to standardize in the area of non-standard items. A method
 3124 that exists on some historical implementations is the use of the so-called `/local/bin` or
 3125 `/usr/local/bin` directory to separate local or additional copies or versions of utilities. Another
 3126 method that is also used is to isolate utilities into completely separate domains. Still another

3127 method to ensure that the desired utility is being used is to request the utility by its full
 3128 pathname. There are many approaches to this situation; the examples given above serve to
 3129 illustrate that there is more than one.

3130 **A.13 Headers**

3131 **A.13.1 Format of Entries**

3132 Each header reference page has a common layout of sections describing the interface. This
 3133 layout is similar to the manual page or “man” page format shipped with most UNIX systems,
 3134 and each header has sections describing the SYNOPSIS and DESCRIPTION. These are the two
 3135 sections that relate to conformance.

3136 Additional sections are informative, and add considerable information for the application
 3137 developer. APPLICATION USAGE sections provide additional caveats, issues, and
 3138 recommendations to the developer. RATIONALE sections give additional information on the
 3139 decisions made in defining the interface.

3140 FUTURE DIRECTIONS sections act as pointers to related work that may impact the interface in
 3141 the future, and often cautions the developer to architect the code to account for a change in this
 3142 area. Note that a future directions statement should not be taken as a commitment to adopt a
 3143 feature or interface in the future.

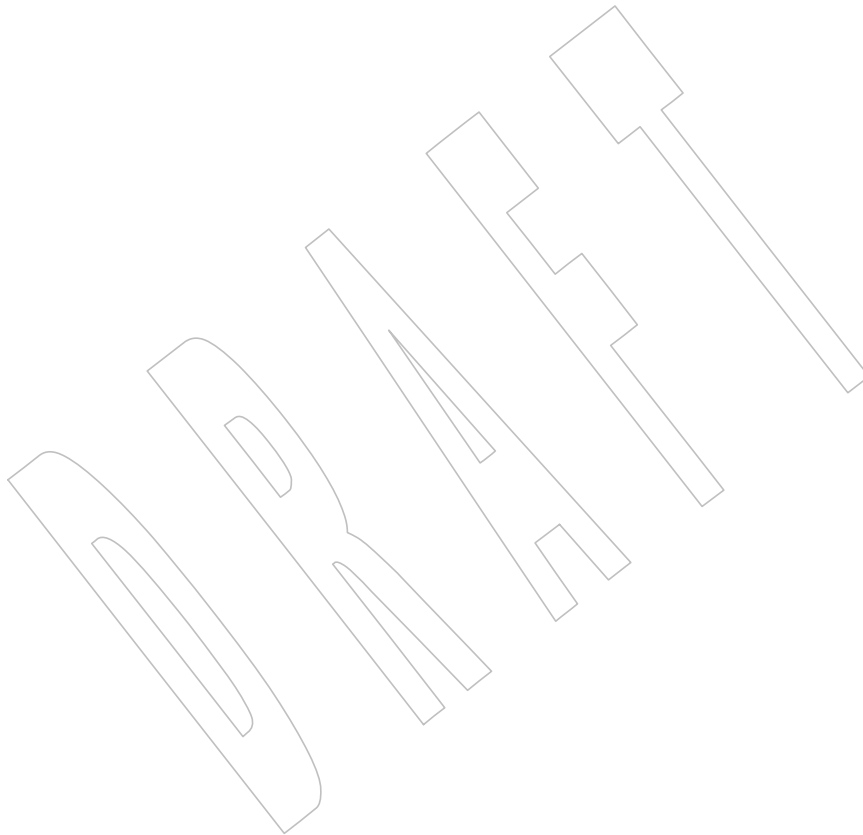
3144 The CHANGE HISTORY section describes when the interface was introduced, and how it has
 3145 changed.

3146 Option labels and margin markings in the page can be useful in guiding the application
 3147 developer.

3148 **A.13.2 Removed Headers in Issue 7**

3149 The headers removed in Issue 7 (from the Issue 6 base document) are as follows:

3150 Removed Headers in Issue 7	
3151 <sys/timeb.h>	<ucontext.h>



3152

Rationale (Informative)

3153

Part B:

3154

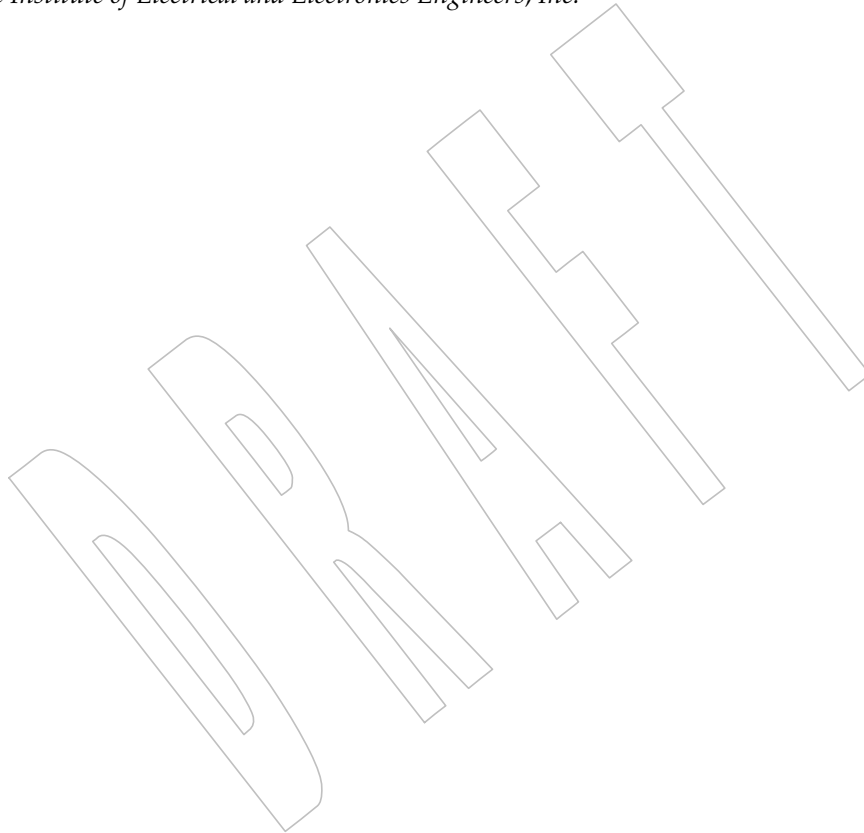
System Interfaces

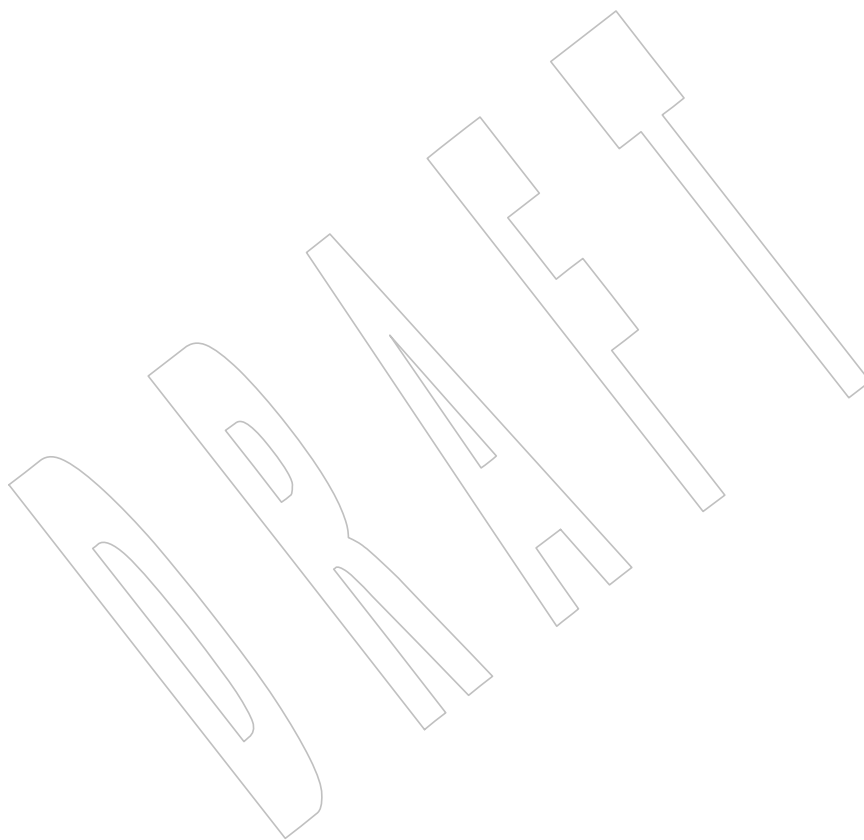
3155

The Open Group

3156

The Institute of Electrical and Electronics Engineers, Inc.





*Rationale for System Interfaces***B.1 Introduction****B.1.1 Scope**

Refer to [Section A.1.1](#) (on page 3).

B.1.2 Conformance

Refer to [Section A.2](#) (on page 9).

B.1.3 Normative References

There is no additional rationale provided for this section.

B.1.4 Change History

The change history is provided as an informative section, to track changes from previous issues of IEEE Std 1003.1-200x.

The following sections describe changes made to the System Interfaces volume of IEEE Std 1003.1-200x since Issue 6 of the base document. The CHANGE HISTORY section for each entry details the technical changes that have been made to that entry from Issue 5. Changes between earlier issues of the base document and Issue 5 are not included.

Changes from Issue 6 to Issue 7 (IEEE Std 1003.1-200x)

The following list summarizes the major changes that were made in the System Interfaces volume of IEEE Std 1003.1-200x from Issue 6 to Issue 7:

- The Open Group Technical Standard, 2006, Extended API Set Part 1 is incorporated.
- The Open Group Technical Standard, 2006, Extended API Set Part 2 is incorporated.
- The Open Group Technical Standard, 2006, Extended API Set Part 3 is incorporated.
- The Open Group Technical Standard, 2006, Extended API Set Part 4 is incorporated.
- Existing functionality is aligned with ISO/IEC 9899:1999, Programming Languages — C, ISO/IEC 9899:1999/Cor.2:2004(E)
- Austin Group defect reports, IEEE Interpretations against IEEE Std 1003.1, and responses to ISO/IEC defect reports against ISO/IEC 9945 are applied.
- The Open Group corrigenda and resolutions are applied.
- Features, marked legacy or obsolescent in the base document, have been considered for withdrawal in the revision.
- The options within the standard have been revised.

3188

New Features in Issue 7

3189

The functions first introduced in Issue 7 (over the Issue 6 base document) are as follows:

3190

3191

3192

3193

3194

3195

3196

3197

3198

3199

3200

3201

3202

3203

3204

3205

3206

3207

3208

3209

3210

3211

3212

3213

3214

3215

3216

3217

3218

3219

3220

New Functions in Issue 7		
<i>alphasort()</i>	<i>iswctype_l()</i>	<i>strcoll_l()</i>
<i>dirfd()</i>	<i>iswdigit_l()</i>	<i>strfmon_l()</i>
<i>dprintf()</i>	<i>iswgraph_l()</i>	<i>strncasecmp_l()</i>
<i>duplocale()</i>	<i>iswlower_l()</i>	<i>strndup()</i>
<i>faccessat()</i>	<i>iswprint_l()</i>	<i>strlen()</i>
<i>fchmodat()</i>	<i>iswpunct_l()</i>	<i>strsignal()</i>
<i>fchownat()</i>	<i>iswspace_l()</i>	<i>strxfrm_l()</i>
<i>fdopendir()</i>	<i>iswupper_l()</i>	<i>symlinkat()</i>
<i>fexecve()</i>	<i>iswxdigit_l()</i>	<i>tolower_l()</i>
<i>fnemopen()</i>	<i>isxdigit_l()</i>	<i>toupper_l()</i>
<i>freelocale()</i>	<i>linkat()</i>	<i>towctrans_l()</i>
<i>fstatat()</i>	<i>mbsnrtowcs()</i>	<i>towlower()</i>
<i>futimesat()</i>	<i>mkdirat()</i>	<i>towupper()</i>
<i>getdelim()</i>	<i>mkdtemp()</i>	<i>unlinkat()</i>
<i>getline()</i>	<i>mkffloat()</i>	<i>uselocale()</i>
<i>isalnum_l()</i>	<i>mknodat()</i>	<i>wcpcpy()</i>
<i>isalpha_l()</i>	<i>newlocale()</i>	<i>wcpncpy()</i>
<i>isblank_l()</i>	<i>openat()</i>	<i>wcscasecmp()</i>
<i>iscntrl_l()</i>	<i>open_memstream()</i>	<i>wcscasecmp_l()</i>
<i>isdigit_l()</i>	<i>psiginfo()</i>	<i>wcscoll_l()</i>
<i>isgraph_l()</i>	<i>psignal()</i>	<i>wcsdup()</i>
<i>islower_l()</i>	<i>pthread_mutexattr_getrobust()</i>	<i>wcsncasecmp()</i>
<i>isprint_l()</i>	<i>pthread_mutexattr_setrobust()</i>	<i>wcsncasecmp_l()</i>
<i>ispunct_l()</i>	<i>pthread_mutex_consistent()</i>	<i>wcsnlen()</i>
<i>isspace_l()</i>	<i>readlinkat()</i>	<i>wcsnrtombs()</i>
<i>isupper_l()</i>	<i>renameat()</i>	<i>wcsxfrm_l()</i>
<i>iswalnum_l()</i>	<i>scandir()</i>	<i>wctrans_l()</i>
<i>iswalphal_l()</i>	<i>stpncpy()</i>	<i>wctype_l()</i>
<i>iswblank_l()</i>	<i>stpncpy()</i>	
<i>iswcntrl_l()</i>	<i>strcasecmp_l()</i>	

3221

Newly Mandated Functions in Issue 7

3222

The functions that were previously part of an option group but are now mandatory in Issue 7 are as follows:

3223

3224

Newly Mandated Functions in Issue 7

3225	<i>aio_cancel()</i>	<i>pthread_atfork()</i>	<i>pthread_rwlock_tryrdlock()</i>
3226	<i>aio_error()</i>	<i>pthread_attr_destroy()</i>	<i>pthread_rwlock_trywrlock()</i>
3227	<i>aio_fsync()</i>	<i>pthread_attr_getdetachstate()</i>	<i>pthread_rwlock_unlock()</i>
3228	<i>aio_read()</i>	<i>pthread_attr_getguardsize()</i>	<i>pthread_rwlock_wrlock()</i>
3229	<i>aio_return()</i>	<i>pthread_attr_getschedparam()</i>	<i>pthread_rwlockattr_destroy()</i>
3230	<i>aio_suspend()</i>	<i>pthread_attr_init()</i>	<i>pthread_rwlockattr_init()</i>
3231	<i>aio_write()</i>	<i>pthread_attr_setdetachstate()</i>	<i>pthread_self()</i>
3232	<i>asctime_r()</i>	<i>pthread_attr_setguardsize()</i>	<i>pthread_setcancelstate()</i>
3233	<i>catclose()</i>	<i>pthread_attr_setschedparam()</i>	<i>pthread_setcanceltype()</i>
3234	<i>catgets()</i>	<i>pthread_barrier_destroy()</i>	<i>pthread_setspecific()</i>
3235	<i>catopen()</i>	<i>pthread_barrier_init()</i>	<i>pthread_spin_destroy()</i>
3236	<i>clock_getres()</i>	<i>pthread_barrier_wait()</i>	<i>pthread_spin_init()</i>
3237	<i>clock_gettime()</i>	<i>pthread_barrierattr_destroy()</i>	<i>pthread_spin_lock()</i>
3238	<i>clock_nanosleep()</i>	<i>pthread_barrierattr_init()</i>	<i>pthread_spin_trylock()</i>
3239	<i>clock_settime()</i>	<i>pthread_cancel()</i>	<i>pthread_spin_unlock()</i>
3240	<i>ctime_r()</i>	<i>pthread_cleanup_pop()</i>	<i>pthread_testcancel()</i>
3241	<i>dlclose()</i>	<i>pthread_cleanup_push()</i>	<i>putc_unlocked()</i>
3242	<i>dlderror()</i>	<i>pthread_cond_broadcast()</i>	<i>putchar_unlocked()</i>
3243	<i>dlopen()</i>	<i>pthread_cond_destroy()</i>	<i>pwrite()</i>
3244	<i>dlsym()</i>	<i>pthread_cond_init()</i>	<i>rand_r()</i>
3245	<i>fehdir()</i>	<i>pthread_cond_signal()</i>	<i>readdir_r()</i>
3246	<i>flockfile()</i>	<i>pthread_cond_timedwait()</i>	<i>sem_close()</i>
3247	<i>fstatvfs()</i>	<i>pthread_cond_wait()</i>	<i>sem_destroy()</i>
3248	<i>ftrylockfile()</i>	<i>pthread_condattr_destroy()</i>	<i>sem_getvalue()</i>
3249	<i>funlockfile()</i>	<i>pthread_condattr_getclock()</i>	<i>sem_init()</i>
3250	<i>getc_unlocked()</i>	<i>pthread_condattr_init()</i>	<i>sem_open()</i>
3251	<i>getchar_unlocked()</i>	<i>pthread_condattr_setclock()</i>	<i>sem_post()</i>
3252	<i>getgrgid_r()</i>	<i>pthread_create()</i>	<i>sem_timedwait()</i>
3253	<i>getgrnam_r()</i>	<i>pthread_detach()</i>	<i>sem_trywait()</i>
3254	<i>getlogin_r()</i>	<i>pthread_equal()</i>	<i>sem_unlink()</i>
3255	<i>getpgid()</i>	<i>pthread_exit()</i>	<i>sem_wait()</i>
3256	<i>getpwnam_r()</i>	<i>pthread_getspecific()</i>	<i>sigqueue()</i>
3257	<i>getpwuid_r()</i>	<i>pthread_join()</i>	<i>sigqueue()</i>
3258	<i>getsid()</i>	<i>pthread_key_create()</i>	<i>sigtimedwait()</i>
3259	<i>getsubopt()</i>	<i>pthread_key_delete()</i>	<i>sigwaitinfo()</i>
3260	<i>gmtime_r()</i>	<i>pthread_mutex_destroy()</i>	<i>statvfs()</i>
3261	<i>iconv()</i>	<i>pthread_mutex_init()</i>	<i>strcasemp()</i>
3262	<i>iconv_close()</i>	<i>pthread_mutex_lock()</i>	<i>strdup()</i>
3263	<i>iconv_open()</i>	<i>pthread_mutex_timedlock()</i>	<i>strerror_r()</i>
3264	<i>lchown()</i>	<i>pthread_mutex_trylock()</i>	<i>strfnon()</i>
3265	<i>lio_listio()</i>	<i>pthread_mutex_unlock()</i>	<i>strncasemp()</i>
3266	<i>localtime_r()</i>	<i>pthread_mutexattr_destroy()</i>	<i>strtok_r()</i>
3267	<i>mkstemp()</i>	<i>pthread_mutexattr_gettype()</i>	<i>tcgetsid()</i>
3268	<i>mmap()</i>	<i>pthread_mutexattr_init()</i>	<i>timer_create()</i>
3269	<i>mprotect()</i>	<i>pthread_mutexattr_settype()</i>	<i>timer_delete()</i>
3270	<i>munmap()</i>	<i>pthread_once()</i>	<i>timer_getoverrun()</i>
3271	<i>nanosleep()</i>	<i>pthread_rwlock_destroy()</i>	<i>timer_gettime()</i>
3272	<i>nl_langinfo()</i>	<i>pthread_rwlock_init()</i>	<i>timer_settime()</i>
3273	<i>poll()</i>	<i>pthread_rwlock_rdlock()</i>	<i>truncate()</i>
3274	<i>posix_trace_timedgetnext_event()</i>	<i>pthread_rwlock_timedrdlock()</i>	<i>ttynam_r()</i>
3275	<i>pread()</i>	<i>pthread_rwlock_timedwrlock()</i>	<i>waitid()</i>

3276

Obsolescent Functions in Issue 7

3277

The base functions moved to obsolescent status in Issue 7 (from the Issue 6 base document) are as follows:

3278

3279

Obsolescent Base Functions in Issue 7

3280

asctime() *ctime_r()*

3281

asctime_r() *tmpnam()*

3282

ctime()

3283

The XSI functions moved to obsolescent status in Issue 7 (from the Issue 6 base document) are as follows:

3284

3285

Obsolescent XSI Functions in Issue 7

3286

_longjmp() *setpgrp()*

3287

_setjmp() *sighold()*

3288

_tolower() *sigignore()*

3289

_toupper() *sigpause()*

3290

ftw() *sigrelse()*

3291

getitimer() *sigset()*

3292

gettimeofday() *siginterrupt()*

3293

isascii() *tempnam()*

3294

pthread_getconcurrency() *toascii()*

3295

pthread_setconcurrency() *ulimit()*

3296

setitimer()

3297

Removed Functions and Symbols in Issue 7

3298

The functions and symbols removed in Issue 7 (from the Issue 6 base document) are as follows:

3299

Removed Functions and Symbols in Issue 7

3300

bcmp() *gethostbyaddr()* *rindex()*

3301

bcopy() *gethostbyname()* *scalb()*

3302

bsd_signal() *getwd()* *setcontext()*

3303

bzero() *h_errno* *swapcontext()*

3304

ecvt() *index()* *ualarm()*

3305

fcvt() *makecontext()* *usleep()*

3306

ftime() *mktemp()* *vfork()*

3307

gcvt() *pthread_attr_getstackaddr()* *wcswcs()*

3308

getcontext() *pthread_attr_setstackaddr()*

3309

B.1.5 Terminology

3310

Refer to [Section A.1.4](#) (on page 6).

3311 **B.1.6 Definitions**

3312 Refer to [Section A.3](#) (on page 13).

3313 **B.1.7 Relationship to Other Formal Standards**

3314 There is no additional rationale provided for this section.

3315 **B.1.8 Portability**

3316 Refer to [Section A.1.5](#) (on page 8).

3317 **B.1.8.1 Codes**

3318 Refer to [Section A.1.5.1](#) (on page 8).

3319 **B.1.9 Format of Entries**

3320 Each system interface reference page has a common layout of sections describing the interface.
 3321 This layout is similar to the manual page or “man” page format shipped with most UNIX
 3322 systems, and each header has sections describing the SYNOPSIS, DESCRIPTION, RETURN
 3323 VALUE, and ERRORS. These are the four sections that relate to conformance.

3324 Additional sections are informative, and add considerable information for the application
 3325 developer. EXAMPLES sections provide example usage. APPLICATION USAGE sections
 3326 provide additional caveats, issues, and recommendations to the developer. RATIONALE
 3327 sections give additional information on the decisions made in defining the interface.

3328 FUTURE DIRECTIONS sections act as pointers to related work that may impact the interface in
 3329 the future, and often cautions the developer to architect the code to account for a change in this
 3330 area. Note that a future directions statement should not be taken as a commitment to adopt a
 3331 feature or interface in the future.

3332 The CHANGE HISTORY section describes when the interface was introduced, and how it has
 3333 changed.

3334 Option labels and margin markings in the page can be useful in guiding the application
 3335 developer.

3336 **B.2 General Information**

3337 **B.2.1 Use and Implementation of Functions**

3338 The information concerning the use of functions was adapted from a description in the ISO C
 3339 standard. Here is an example of how an application program can protect itself from functions
 3340 that may or may not be macros, rather than true functions:

3341 The `atoi()` function may be used in any of several ways:

- 3342 • By use of its associated header (possibly generating a macro expansion):

```
3343 #include <stdlib.h>
3344 /* ... */
3345 i = atoi(str);
```

- 3346 • By use of its associated header (assuredly generating a true function call):

```
3347 #include <stdlib.h>
3348 #undef atoi
```



```
3349     /* ... */
3350     i = atoi(str);
```

```
3351     or:
```

```
3352     #include <stdlib.h>
3353     /* ... */
3354     i = (atoi) (str);
```

```
3355     • By explicit declaration:
```

```
3356     extern int atoi (const char *);
3357     /* ... */
3358     i = atoi(str);
```

```
3359     • By implicit declaration:
```

```
3360     /* ... */
3361     i = atoi(str);
```

```
3362     (Assuming no function prototype is in scope. This is not allowed by the ISO C standard for
3363     functions with variable arguments; furthermore, parameter type conversion “widening” is
3364     subject to different rules in this case.)
```

```
3365     Note that the ISO C standard reserves names starting with ‘_’ for the compiler. Therefore, the
3366     compiler could, for example, implement an intrinsic, built-in function _asm_builtin_atoi(), which
3367     it recognized and expanded into inline assembly code. Then, in <stdlib.h>, there could be the
3368     following:
```

```
3369     #define atoi(X) _asm_builtin_atoi(X)
```

```
3370     The user’s “normal” call to atoi() would then be expanded inline, but the implementor would
3371     also be required to provide a callable function named atoi() for use when the application
3372     requires it; for example, if its address is to be stored in a function pointer variable.
```

3373 B.2.2 The Compilation Environment

3374 B.2.2.1 POSIX.1 Symbols

```
3375     This and the following section address the issue of “name space pollution”. The ISO C standard
3376     requires that the name space beyond what it reserves not be altered except by explicit action of
3377     the application writer. This section defines the actions to add the POSIX.1 symbols for those
3378     headers where both the ISO C standard and POSIX.1 need to define symbols, and also where the
3379     XSI option extends the base standard.
```

```
3380     When headers are used to provide symbols, there is a potential for introducing symbols that the
3381     application writer cannot predict. Ideally, each header should only contain one set of symbols,
3382     but this is not practical for historical reasons. Thus, the concept of feature test macros is
3383     included. Two feature test macros are explicitly defined by IEEE Std 1003.1-200x; it is expected
3384     that future revisions may add to this.
```

```
3385     Note: Feature test macros allow an application to announce to the implementation its desire to have
3386     certain symbols and prototypes exposed. They should not be confused with the version test
3387     macros and constants for options in <unistd.h> which are the implementation’s way of
3388     announcing functionality to the application.
```

```
3389     It is further intended that these feature test macros apply only to the headers specified by
3390     IEEE Std 1003.1-200x. Implementations are expressly permitted to make visible symbols not
3391     specified by IEEE Std 1003.1-200x, within both POSIX.1 and other headers, under the control of
3392     feature test macros that are not defined by IEEE Std 1003.1-200x.
```

3393 **The `_POSIX_C_SOURCE` Feature Test Macro**

3394 Since `_POSIX_SOURCE` specified by the POSIX.1-1990 standard did not have a value associated
 3395 with it, the `_POSIX_C_SOURCE` macro replaces it, allowing an application to inform the system
 3396 of the revision of the standard to which it conforms. This symbol will allow implementations to
 3397 support various revisions of IEEE Std 1003.1-200x simultaneously. For instance, when either
 3398 `_POSIX_SOURCE` is defined or `_POSIX_C_SOURCE` is defined as 1, the system should make
 3399 visible the same name space as permitted and required by the POSIX.1-1990 standard. When
 3400 `_POSIX_C_SOURCE` is defined, the state of `_POSIX_SOURCE` is completely irrelevant.

3401 It is expected that C bindings to future POSIX standards will define new values for
 3402 `_POSIX_C_SOURCE`, with each new value reserving the name space for that new standard, plus
 3403 all earlier POSIX standards.

3404 **The `_XOPEN_SOURCE` Feature Test Macro**

3405 The feature test macro `_XOPEN_SOURCE` is provided as the announcement mechanism for the
 3406 application that it requires functionality from the Single UNIX Specification. `_XOPEN_SOURCE`
 3407 must be defined to the value 700 before the inclusion of any header to enable the functionality in
 3408 the Single UNIX Specification. Its definition subsumes the use of `_POSIX_SOURCE` and
 3409 `_POSIX_C_SOURCE`.

3410 An extract of code from a conforming application, that appears before any **#include** statements,
 3411 is given below:

```
3412 #define _XOPEN_SOURCE 700 /* Single UNIX Specification, Version x */
3413 #include ...
```

3414 Note that the definition of `_XOPEN_SOURCE` with the value 700 makes the definition of
 3415 `_POSIX_C_SOURCE` redundant and it can safely be omitted.

3416 **B.2.2.2 The Name Space**

3417 The reservation of identifiers is paraphrased from the ISO C standard. The text is included
 3418 because it needs to be part of IEEE Std 1003.1-200x, regardless of possible changes in future
 3419 versions of the ISO C standard.

3420 These identifiers may be used by implementations, particularly for feature test macros.
 3421 Implementations should not use feature test macro names that might be reasonably used by a
 3422 standard.

3423 Including headers more than once is a reasonably common practice, and it should be carried
 3424 forward from the ISO C standard. More significantly, having definitions in more than one
 3425 header is explicitly permitted. Where the potential declaration is “benign” (the same definition
 3426 twice) the declaration can be repeated, if that is permitted by the compiler. (This is usually true
 3427 of macros, for example.) In those situations where a repetition is not benign (for example,
 3428 **typedefs**), conditional compilation must be used. The situation actually occurs both within the
 3429 ISO C standard and within POSIX.1: **time_t** should be in **<sys/types.h>**, and the ISO C standard
 3430 mandates that it be in **<time.h>**.

3431 The area of name space pollution *versus* additions to structures is difficult because of the macro
 3432 structure of C. The following discussion summarizes all the various problems with and
 3433 objections to the issue.

3434 Note the phrase “user-defined macro”. Users are not permitted to define macro names (or any
 3435 other name) beginning with “_**[A-Z]**”. Thus, the conflict cannot occur for symbols reserved
 3436 to the vendor’s name space, and the permission to add fields automatically applies, without
 3437 qualification, to those symbols.

- 3438 1. Data structures (and unions) need to be defined in headers by implementations to meet
3439 certain requirements of POSIX.1 and the ISO C standard.
- 3440 2. The structures defined by POSIX.1 are typically minimal, and any practical
3441 implementation would wish to add fields to these structures either to hold additional
3442 related information or for backwards-compatibility (or both). Future standards (and *de*
3443 *facto* standards) would also wish to add to these structures. Issues of field alignment
3444 make it impractical (at least in the general case) to simply omit fields when they are not
3445 defined by the particular standard involved.

3446 The **dirent** structure is an example of such a minimal structure (although one could argue
3447 about whether the other fields need visible names). The *st_rdev* field of most
3448 implementations' **stat** structure is a common example where extension is needed and
3449 where a conflict could occur.

- 3450 3. Fields in structures are in an independent name space, so the addition of such fields
3451 presents no problem to the C language itself in that such names cannot interact with
3452 identically named user symbols because access is qualified by the specific structure name.
- 3453 4. There is an exception to this: macro processing is done at a lexical level. Thus, symbols
3454 added to a structure might be recognized as user-provided macro names at the location
3455 where the structure is declared. This only can occur if the user-provided name is declared
3456 as a macro before the header declaring the structure is included. The user's use of the
3457 name after the declaration cannot interfere with the structure because the symbol is
3458 hidden and only accessible through access to the structure. Presumably, the user would
3459 not declare such a macro if there was an intention to use that field name.
- 3460 5. Macros from the same or a related header might use the additional fields in the structure,
3461 and those field names might also collide with user macros. Although this is a less
3462 frequent occurrence, since macros are expanded at the point of use, no constraint on the
3463 order of use of names can apply.
- 3464 6. An "obvious" solution of using names in the reserved name space and then redefining
3465 them as macros when they should be visible does not work because this has the effect of
3466 exporting the symbol into the general name space. For example, given a (hypothetical)
3467 system-provided header **<h.h>**, and two parts of a C program in **a.c** and **b.c**, in header
3468 **<h.h>**:

```
3469 struct foo {
3470     int __i;
3471 }
3472 #ifdef _FEATURE_TEST
3473 #define i __i;
3474 #endif
```

3475 In file **a.c**:

```
3476 #include h.h
3477 extern int i;
3478 ...
```

3479 In file **b.c**:

```
3480 extern int i;
3481 ...
```

3482 The symbol that the user thinks of as *i* in both files has an external name of `__i` in **a.c**; the
3483 same symbol *i* in **b.c** has an external name *i* (ignoring any hidden manipulations the

3484 compiler might perform on the names). This would cause a mysterious name resolution
3485 problem when **a.o** and **b.o** are linked.

3486 Simply avoiding definition then causes alignment problems in the structure.

3487 A structure of the form:

```
3488 struct foo {
3489     union {
3490         int __i;
3491 #ifdef _FEATURE_TEST
3492         int i;
3493 #endif
3494     } __ii;
3495 }
```

3496 does not work because the name of the logical field *i* is *__ii.i*, and introduction of a macro
3497 to restore the logical name immediately reintroduces the problem discussed previously
3498 (although its manifestation might be more immediate because a syntax error would result
3499 if a recursive macro did not cause it to fail first).

3500 7. A more workable solution would be to declare the structure:

```
3501 struct foo {
3502 #ifdef _FEATURE_TEST
3503     int i;
3504 #else
3505     int __i;
3506 #endif
3507 }
```

3508 However, if a macro (particularly one required by a standard) is to be defined that uses
3509 this field, two must be defined: one that uses *i*, the other that uses *__i*. If more than one
3510 additional field is used in a macro and they are conditional on distinct combinations of
3511 features, the complexity goes up as 2^n .

3512 All this leaves a difficult situation: vendors must provide very complex headers to deal with
3513 what is conceptually simple and safe—adding a field to a structure. It is the possibility of user-
3514 provided macros with the same name that makes this difficult.

3515 Several alternatives were proposed that involved constraining the user's access to part of the
3516 name space available to the user (as specified by the ISO C standard). In some cases, this was
3517 only until all the headers had been included. There were two proposals discussed that failed to
3518 achieve consensus:

- 3519 1. Limiting it for the whole program.
- 3520 2. Restricting the use of identifiers containing only uppercase letters until after all system
3521 headers had been included. It was also pointed out that because macros might wish to
3522 access fields of a structure (and macro expansion occurs totally at point of use) restricting
3523 names in this way would not protect the macro expansion, and thus the solution was
3524 inadequate.

3525 It was finally decided that reservation of symbols would occur, but as constrained.

3526 The current wording also allows the addition of fields to a structure, but requires that user
3527 macros of the same name not interfere. This allows vendors to do one of the following:

- 3528 • Not create the situation (do not extend the structures with user-accessible names or use the
- 3529 solution in (7) above)
- 3530 • Extend their compilers to allow some way of adding names to structures and macros safely

3531 There are at least two ways that the compiler might be extended: add new preprocessor
 3532 directives that turn off and on macro expansion for certain symbols (without changing the value
 3533 of the macro) and a function or lexical operation that suppresses expansion of a word. The latter
 3534 seems more flexible, particularly because it addresses the problem in macros as well as in
 3535 declarations.

3536 The following seems to be a possible implementation extension to the C language that will do
 3537 this: any token that during macro expansion is found to be preceded by three '#' symbols shall
 3538 not be further expanded in exactly the same way as described for macros that expand to their
 3539 own name as in Section 3.8.3.4 of the ISO C standard. A vendor may also wish to implement this
 3540 as an operation that is lexically a function, which might be implemented as:

```
3541 #define __safe_name(x) ###x
```

3542 Using a function notation would insulate vendors from changes in standards until such a
 3543 functionality is standardized (if ever). Standardization of such a function would be valuable
 3544 because it would then permit third parties to take advantage of it portably in software they may
 3545 supply.

3546 The symbols that are “explicitly permitted, but not required by IEEE Std 1003.1-200x” include
 3547 those classified below. (That is, the symbols classified below might, but are not required to, be
 3548 present when `_POSIX_C_SOURCE` is defined to have the value `200xxxL`.)

- 3549 • Symbols in `<limits.h>` and `<unistd.h>` that are defined to indicate support for options or
- 3550 limits that are constant at compile-time
- 3551 • Symbols in the name space reserved for the implementation by the ISO C standard
- 3552 • Symbols in a name space reserved for a particular type of extension (for example, type
- 3553 names ending with `_t` in `<sys/types.h>`)
- 3554 • Additional members of structures or unions whose names do not reduce the name space
- 3555 reserved for applications

3556 Since both implementations and future revisions of IEEE Std 1003.1 and other POSIX standards
 3557 may use symbols in the reserved spaces described in these tables, there is a potential for name
 3558 space clashes. To avoid future name space clashes when adding symbols, implementations
 3559 should not use the `posix_`, `POSIX_`, or `_POSIX_` prefixes.

3560 IEEE Std 1003.1-2001/Cor 1-2002, item XSH/TC1/D6/2 is applied, deleting the entries `POSIX_`,
 3561 `_POSIX_`, and `posix_` from the column of allowed name space prefixes for use by an
 3562 implementation in the first table. The presence of these prefixes was contradicting later text
 3563 which states that: “The prefixes `posix_`, `POSIX_`, and `_POSIX_` are reserved for use by Shell and
 3564 Utilities volume of IEEE Std 1003.1-200x, Chapter 2, Shell Command Language and other POSIX
 3565 standards. Implementations may add symbols to the headers shown in the following table,
 3566 provided the identifiers ... do not use the reserved prefixes `posix_`, `POSIX_`, or `_POSIX_`.”

3567 IEEE Std 1003.1-2001/Cor 1-2002, item XSH/TC1/D6/3 is applied, correcting the reserved
 3568 macro prefix from: “`PRI[a-z]`, `SCN[a-z]`” to: “`PRI[Xa-z]`, `SCN[Xa-z]`” in the second table. The
 3569 change was needed since the ISO C standard allows implementations to define macros of the
 3570 form `PRI` or `SCN` followed by any lowercase letter or ‘`X`’ in `<inttypes.h>`. (The
 3571 ISO/IEC 9899:1999 standard, Subclause 7.26.4.)

3572 IEEE Std 1003.1-2001/Cor 1-2002, item XSH/TC1/D6/4 is applied, adding a new section listing
 3573 reserved names for the `<stdint.h>` header. This change is for alignment with the ISO C standard.

3574 IEEE Std 1003.1-2001/Cor 2-2004, item XSH/TC2/D6/2 is applied, making it clear that
3575 implementations are permitted to have symbols with the prefix `_POSIX_` visible in any header.

3576 IEEE Std 1003.1-2001/Cor 2-2004, item XSH/TC2/D6/3 is applied, updating the table of
3577 allowed macro prefixes to include the prefix `FP_[A-Z]` for `<math.h>`. This text is added for
3578 consistency with the `<math.h>` reference page in the Base Definitions volume of
3579 IEEE Std 1003.1-200x which permits additional implementation-defined floating-point
3580 classifications.

3581 Austin Group Interpretation 1003.1-2001 #048 is applied, reserving `SEEK_` in the name space.

3582 B.2.3 Error Numbers

3583 It was the consensus of the standard developers that to allow the conformance document to state
3584 that an error occurs and under what conditions, but to disallow a statement that it never occurs,
3585 does not make sense. It could be implied by the current wording that this is allowed, but to
3586 reduce the possibility of future interpretation requests, it is better to make an explicit statement.

3587 The ISO C standard requires that `errno` be an assignable lvalue. Originally, the definition in
3588 POSIX.1 was stricter than that in the ISO C standard, `extern int errno`, in order to support
3589 historical usage. In a multi-threaded environment, implementing `errno` as a global variable
3590 results in non-deterministic results when accessed. It is required, however, that `errno` work as a
3591 per-thread error reporting mechanism. In order to do this, a separate `errno` value has to be
3592 maintained for each thread. The following section discusses the various alternative solutions
3593 that were considered.

3594 In order to avoid this problem altogether for new functions, these functions avoid using `errno`
3595 and, instead, return the error number directly as the function return value; a return value of zero
3596 indicates that no error was detected.

3597 For any function that can return errors, the function return value is not used for any purpose
3598 other than for reporting errors. Even when the output of the function is scalar, it is passed
3599 through a function argument. While it might have been possible to allow some scalar outputs to
3600 be coded as negative function return values and mixed in with positive error status returns, this
3601 was rejected—using the return value for a mixed purpose was judged to be of limited use and
3602 error prone.

3603 Checking the value of `errno` alone is not sufficient to determine the existence or type of an error,
3604 since it is not required that a successful function call clear `errno`. The variable `errno` should only
3605 be examined when the return value of a function indicates that the value of `errno` is meaningful.
3606 In that case, the function is required to set the variable to something other than zero.

3607 The variable `errno` is never set to zero by any function call; to do so would contradict the ISO C
3608 standard.

3609 POSIX.1 requires (in the ERRORS sections of function descriptions) certain error values to be set
3610 in certain conditions because many existing applications depend on them. Some error numbers,
3611 such as `[EFAULT]`, are entirely implementation-defined and are noted as such in their
3612 description in the ERRORS section. This section otherwise allows wide latitude to the
3613 implementation in handling error reporting.

3614 Some of the ERRORS sections in IEEE Std 1003.1-200x have two subsections. The first:

3615 “The function shall fail if:”

3616 could be called the “mandatory” section.

3617 The second:

3618 “The function may fail if:”

3619 could be informally known as the “optional” section.

3620 Attempting to infer the quality of an implementation based on whether it detects optional error
3621 conditions is not useful.

3622 Following each one-word symbolic name for an error, there is a description of the error. The
3623 rationale for some of the symbolic names follows:

3624 [ECANCELED] This spelling was chosen as being more common.

3625 [EFAULT] Most historical implementations do not catch an error and set *errno* when an
3626 invalid address is given to the functions *wait()*, *time()*, or *times()*. Some
3627 implementations cannot reliably detect an invalid address. And most systems
3628 that detect invalid addresses will do so only for a system call, not for a library
3629 routine.

3630 [EFTYPE] This error code was proposed in earlier proposals as “Inappropriate operation
3631 for file type”, meaning that the operation requested is not appropriate for the
3632 file specified in the function call. This code was proposed, although the same
3633 idea was covered by [ENOTTY], because the connotations of the name would
3634 be misleading. It was pointed out that the *fcntl()* function uses the error code
3635 [EINVAL] for this notion, and hence all instances of [EFTYPE] were changed
3636 to this code.

3637 [EINTR] POSIX.1 prohibits conforming implementations from restarting interrupted
3638 system calls of conforming applications unless the SA_RESTART flag is in
3639 effect for the signal. However, it does not require that [EINTR] be returned
3640 when another legitimate value may be substituted; for example, a partial
3641 transfer count when *read()* or *write()* are interrupted. This is only given when
3642 the signal-catching function returns normally as opposed to returns by
3643 mechanisms like *longjmp()* or *siglongjmp()*.

3644 [ELOOP] In specifying conditions under which implementations would generate this
3645 error, the following goals were considered:

- 3646 • To ensure that actual loops are detected, including loops that result from
3647 symbolic links across distributed file systems.
- 3648 • To ensure that during pathname resolution an application can rely on
3649 the ability to follow at least {SYMLOOP_MAX} symbolic links in the
3650 absence of a loop.
- 3651 • To allow implementations to provide the capability of traversing more
3652 than {SYMLOOP_MAX} symbolic links in the absence of a loop.
- 3653 • To allow implementations to detect loops and generate the error prior to
3654 encountering {SYMLOOP_MAX} symbolic links.

3655 [ENAMETOOLONG] When a symbolic link is encountered during pathname resolution, the
3656 contents of that symbolic link are used to create a new pathname. The
3657 standard developers intended to allow, but not require, that implementations
3658 enforce the restriction of {PATH_MAX} on the result of this pathname
3659 substitution.
3660

- [ENOMEM] The term “main memory” is not used in POSIX.1 because it is implementation-defined.
- [ENOTSUP] This error code is to be used when an implementation chooses to implement the required functionality of IEEE Std 1003.1-200x but does not support optional facilities defined by IEEE Std 1003.1-200x. The return of [ENOSYS] is to be taken to indicate that the function of the interface is not supported at all; the function will always fail with this error code.
- [ENOTTY] The symbolic name for this error is derived from a time when device control was done by *ioctl()* and that operation was only permitted on a terminal interface. The term “TTY” is derived from “teletypewriter”, the devices to which this error originally applied.
- [EOVERFLOW] Most of the uses of this error code are related to large file support. Typically, these cases occur on systems which support multiple programming environments with different sizes for **off_t**, but they may also occur in connection with remote file systems.
- In addition, when different programming environments have different widths for types such as **int** and **uid_t**, several functions may encounter a condition where a value in a particular environment is too wide to be represented. In that case, this error should be raised. For example, suppose the currently running process has 64-bit **int**, and file descriptor 9 223 372 036 854 775 807 is open and does not have the close-on-exec flag set. If the process then uses *execl()* to *exec* a file compiled in a programming environment with 32-bit **int**, the call to *execl()* can fail with *errno* set to [EOVERFLOW]. A similar failure can occur with *execl()* if any of the user IDs or any of the group IDs to be assigned to the new process image are out of range for the executed file’s programming environment.
- Note, however, that this condition cannot occur for functions that are explicitly described as always being successful, such as *getpid()*.
- [EPIPE] This condition normally generates the signal SIGPIPE; the error is returned if the signal does not terminate the process.
- [EROFS] In historical implementations, attempting to *unlink()* or *rmdir()* a mount point would generate an [EBUSY] error. An implementation could be envisioned where such an operation could be performed without error. In this case, if *either* the directory entry or the actual data structures reside on a read-only file system, [EROFS] is the appropriate error to generate. (For example, changing the link count of a file on a read-only file system could not be done, as is required by *unlink()*, and thus an error should be reported.)

Three err

- 3706 • Implement *errno* as a per-thread integer variable.
- 3707 • Implement *errno* as a service that can access the per-thread error number.
- 3708 • Change all POSIX.1 calls to accept an extra status argument and avoid setting *errno*.
- 3709 • Change all POSIX.1 calls to raise a language exception.

3710 The first option offers the highest level of compatibility with existing practice but requires
 3711 special support in the linker, compiler, and/or virtual memory system to support the new
 3712 concept of thread private variables. When compared with current practice, the third and fourth
 3713 options are much cleaner, more efficient, and encourage a more robust programming style, but
 3714 they require new versions of all of the POSIX.1 functions that might detect an error. The second
 3715 option offers compatibility with existing code that uses the `<errno.h>` header to define the
 3716 symbol *errno*. In this option, *errno* may be a macro defined:

```
3717 #define errno (*__errno())
3718 extern int      *__errno();
```

3719 This option may be implemented as a per-thread variable whereby an *errno* field is allocated in
 3720 the user space object representing a thread, and whereby the function `__errno()` makes a system
 3721 call to determine the location of its user space object and returns the address of the *errno* field of
 3722 that object. Another implementation, one that avoids calling the kernel, involves allocating
 3723 stacks in chunks. The stack allocator keeps a side table indexed by chunk number containing a
 3724 pointer to the thread object that uses that chunk. The `__errno()` function then looks at the stack
 3725 pointer, determines the chunk number, and uses that as an index into the chunk table to find its
 3726 thread object and thus its private value of *errno*. On most architectures, this can be done in four
 3727 to five instructions. Some compilers may wish to implement `__errno()` inline to improve
 3728 performance.

3729 **Disallowing Return of the [EINTR] Error Code**

3730 Many blocking interfaces defined by IEEE Std 1003.1-200x may return [EINTR] if interrupted
 3731 during their execution by a signal handler. Blocking interfaces introduced under the threads
 3732 functionality do not have this property. Instead, they require that the interface appear to be
 3733 atomic with respect to interruption. In particular, clients of blocking interfaces need not handle
 3734 any possible [EINTR] return as a special case since it will never occur. If it is necessary to restart
 3735 operations or complete incomplete operations following the execution of a signal handler, this is
 3736 handled by the implementation, rather than by the application.

3737 Requiring applications to handle [EINTR] errors on blocking interfaces has been shown to be a
 3738 frequent source of often unreproducible bugs, and it adds no compelling value to the available
 3739 functionality. Thus, blocking interfaces introduced for use by multi-threaded programs do not
 3740 use this paradigm. In particular, in none of the functions `flockfile()`, `pthread_cond_timedwait()`,
 3741 `pthread_cond_wait()`, `pthread_join()`, `pthread_mutex_lock()`, and `sigwait()` did providing [EINTR]
 3742 returns add value, or even particularly make sense. Thus, these functions do not provide for an
 3743 [EINTR] return, even when interrupted by a signal handler. The same arguments can be applied
 3744 to `sem_wait()`, `sem_trywait()`, `sigwaitinfo()`, and `sigtimedwait()`, but implementations are
 3745 permitted to return [EINTR] error codes for these functions for compatibility with earlier
 3746 versions of IEEE Std 1003.1. Applications cannot rely on calls to these functions returning
 3747 [EINTR] error codes when signals are delivered to the calling thread, but they should allow for
 3748 the possibility.

3749 Austin Group Interpretation 1003.1-2001 #050 is applied, allowing [ENOTSUP] and
 3750 [EOPNOTSUPP] to be the same values.

3751 B.2.3.1 Additional Error Numbers

3752 The ISO C standard defines the name space for implementations to add additional error
3753 numbers.

3754 B.2.4 Signal Concepts

3755 Historical implementations of signals, using the *signal()* function, have shortcomings that make
3756 them unreliable for many application uses. Because of this, a new signal mechanism, based very
3757 closely on the one of 4.2 BSD and 4.3 BSD, was added to POSIX.1.

3758 Signal Names

3759 The restriction on the actual type used for **sigset_t** is intended to guarantee that these objects can
3760 always be assigned, have their address taken, and be passed as parameters by value. It is not
3761 intended that this type be a structure including pointers to other data structures, as that could
3762 impact the portability of applications performing such operations. A reasonable implementation
3763 could be a structure containing an array of some integer type.

3764 The signals described in IEEE Std 1003.1-200x must have unique values so that they may be
3765 named as parameters of **case** statements in the body of a C-language **switch** clause. However,
3766 implementation-defined signals may have values that overlap with each other or with signals
3767 specified in IEEE Std 1003.1-200x. An example of this is SIGABRT, which traditionally overlaps
3768 some other signal, such as SIGIOT.

3769 SIGKILL, SIGTERM, SIGUSR1, and SIGUSR2 are ordinarily generated only through the explicit
3770 use of the *kill()* function, although some implementations generate SIGKILL under
3771 extraordinary circumstances. SIGTERM is traditionally the default signal sent by the *kill*
3772 command.

3773 The signals SIGBUS, SIGEMT, SIGIOT, SIGTRAP, and SIGSYS were omitted from POSIX.1
3774 because their behavior is implementation-defined and could not be adequately categorized.
3775 Conforming implementations may deliver these signals, but must document the circumstances
3776 under which they are delivered and note any restrictions concerning their delivery. The signals
3777 SIGFPE, SIGILL, and SIGSEGV are similar in that they also generally result only from
3778 programming errors. They were included in POSIX.1 because they do indicate three relatively
3779 well-categorized conditions. They are all defined by the ISO C standard and thus would have to
3780 be defined by any system with an ISO C standard binding, even if not explicitly included in
3781 POSIX.1.

3782 There is very little that a Conforming POSIX.1 Application can do by catching, ignoring, or
3783 masking any of the signals SIGILL, SIGTRAP, SIGIOT, SIGEMT, SIGBUS, SIGSEGV, SIGSYS, or
3784 SIGFPE. They will generally be generated by the system only in cases of programming errors.
3785 While it may be desirable for some robust code (for example, a library routine) to be able to
3786 detect and recover from programming errors in other code, these signals are not nearly sufficient
3787 for that purpose. One portable use that does exist for these signals is that a command interpreter
3788 can recognize them as the cause of termination of a process (with *wait()*) and print an
3789 appropriate message. The mnemonic tags for these signals are derived from their PDP-11 origin.

3790 The signals SIGSTOP, SIGTSTP, SIGTTIN, SIGTTOU, and SIGCONT are provided for job control
3791 and are unchanged from 4.2 BSD. The signal SIGCHLD is also typically used by job control
3792 shells to detect children that have terminated or, as in 4.2 BSD, stopped.

3793 Some implementations, including System V, have a signal named SIGCLD, which is similar to
3794 SIGCHLD in 4.2 BSD. POSIX.1 permits implementations to have a single signal with both
3795 names. POSIX.1 carefully specifies ways in which conforming applications can avoid the
3796 semantic differences between the two different implementations. The name SIGCHLD was
3797 chosen for POSIX.1 because most current application usages of it can remain unchanged in

3798 conforming applications. SIGCLD in System V has more cases of semantics that POSIX.1 does
 3799 not specify, and thus applications using it are more likely to require changes in addition to the
 3800 name change.

3801 The signals SIGUSR1 and SIGUSR2 are commonly used by applications for notification of
 3802 exceptional behavior and are described as “reserved as application-defined” so that such use is
 3803 not prohibited. Implementations should not generate SIGUSR1 or SIGUSR2, except when
 3804 explicitly requested by *kill()*. It is recommended that libraries not use these two signals, as such
 3805 use in libraries could interfere with their use by applications calling the libraries. If such use is
 3806 unavoidable, it should be documented. It is prudent for non-portable libraries to use non-
 3807 standard signals to avoid conflicts with use of standard signals by portable libraries.

3808 There is no portable way for an application to catch or ignore non-standard signals. Some
 3809 implementations define the range of signal numbers, so applications can install signal-catching
 3810 functions for all of them. Unfortunately, implementation-defined signals often cause problems
 3811 when caught or ignored by applications that do not understand the reason for the signal. While
 3812 the desire exists for an application to be more robust by handling all possible signals (even those
 3813 only generated by *kill()*), no existing mechanism was found to be sufficiently portable to include
 3814 in POSIX.1. The value of such a mechanism, if included, would be diminished given that
 3815 SIGKILL would still not be catchable.

3816 A number of new signal numbers are reserved for applications because the two user signals
 3817 defined by POSIX.1 are insufficient for many realtime applications. A range of signal numbers is
 3818 specified, rather than an enumeration of additional reserved signal names, because different
 3819 applications and application profiles will require a different number of application signals. It is
 3820 not desirable to burden all application domains and therefore all implementations with the
 3821 maximum number of signals required by all possible applications. Note that in this context,
 3822 signal numbers are essentially different signal priorities.

3823 The relatively small number of required additional signals, `{_POSIX_RTSIG_MAX}`, was chosen
 3824 so as not to require an unreasonably large signal mask/set. While this number of signals
 3825 defined in POSIX.1 will fit in a single 32-bit word signal mask, it is recognized that most existing
 3826 implementations define many more signals than are specified in POSIX.1 and, in fact, many
 3827 implementations have already exceeded 32 signals (including the “null signal”). Support of
 3828 `{_POSIX_RTSIG_MAX}` additional signals may push some implementation over the single 32-bit
 3829 word line, but is unlikely to push any implementations that are already over that line beyond
 3830 the 64-signal line.

3831 B.2.4.1 *Signal Generation and Delivery*

3832 The terms defined in this section are not used consistently in documentation of historical
 3833 systems. Each signal can be considered to have a lifetime beginning with generation and ending
 3834 with delivery or acceptance. The POSIX.1 definition of “delivery” does not exclude ignored
 3835 signals; this is considered a more consistent definition. This revised text in several parts of
 3836 IEEE Std 1003.1-200x clarifies the distinct semantics of asynchronous signal delivery and
 3837 synchronous signal acceptance. The previous wording attempted to categorize both under the
 3838 term “delivery”, which led to conflicts over whether the effects of asynchronous signal delivery
 3839 applied to synchronous signal acceptance.

3840 Signals generated for a process are delivered to only one thread. Thus, if more than one thread
 3841 is eligible to receive a signal, one has to be chosen. The choice of threads is left entirely up to the
 3842 implementation both to allow the widest possible range of conforming implementations and to
 3843 give implementations the freedom to deliver the signal to the “easiest possible” thread should
 3844 there be differences in ease of delivery between different threads.

3845 Note that should multiple delivery among cooperating threads be required by an application,
 3846 this can be trivially constructed out of the provided single-delivery semantics. The construction

3847 of a *sigwait_multiple()* function that accomplishes this goal is presented with the rationale for
3848 *sigwaitinfo()*.

3849 Implementations should deliver unblocked signals as soon after they are generated as possible.
3850 However, it is difficult for POSIX.1 to make specific requirements about this, beyond those in
3851 *kill()* and *sigprocmask()*. Even on systems with prompt delivery, scheduling of higher priority
3852 processes is always likely to cause delays.

3853 In general, the interval between the generation and delivery of unblocked signals cannot be
3854 detected by an application. Thus, references to pending signals generally apply to blocked,
3855 pending signals. An implementation registers a signal as pending on the process when no
3856 thread has the signal unblocked and there are no threads blocked in a *sigwait()* function for that
3857 signal. Thereafter, the implementation delivers the signal to the first thread that unblocks the
3858 signal or calls a *sigwait()* function on a signal set containing this signal rather than choosing the
3859 recipient thread at the time the signal is sent.

3860 In the 4.3 BSD system, signals that are blocked and set to SIG_IGN are discarded immediately
3861 upon generation. For a signal that is ignored as its default action, if the action is SIG_DFL and
3862 the signal is blocked, a generated signal remains pending. In the 4.1 BSD system and in
3863 System V Release 3 (two other implementations that support a somewhat similar signal
3864 mechanism), all ignored blocked signals remain pending if generated. Because it is not normally
3865 useful for an application to simultaneously ignore and block the same signal, it was unnecessary
3866 for POSIX.1 to specify behavior that would invalidate any of the historical implementations.

3867 There is one case in some historical implementations where an unblocked, pending signal does
3868 not remain pending until it is delivered. In the System V implementation of *signal()*, pending
3869 signals are discarded when the action is set to SIG_DFL or a signal-catching routine (as well as
3870 to SIG_IGN). Except in the case of setting SIGCHLD to SIG_DFL, implementations that do this
3871 do not conform completely to POSIX.1. Some earlier proposals for POSIX.1 explicitly stated this,
3872 but these statements were redundant due to the requirement that functions defined by POSIX.1
3873 not change attributes of processes defined by POSIX.1 except as explicitly stated.

3874 POSIX.1 specifically states that the order in which multiple, simultaneously pending signals are
3875 delivered is unspecified. This order has not been explicitly specified in historical
3876 implementations, but has remained quite consistent and been known to those familiar with the
3877 implementations. Thus, there have been cases where applications (usually system utilities) have
3878 been written with explicit or implicit dependencies on this order. Implementors and others
3879 porting existing applications may need to be aware of such dependencies.

3880 When there are multiple pending signals that are not blocked, implementations should arrange
3881 for the delivery of all signals at once, if possible. Some implementations stack calls to all pending
3882 signal-catching routines, making it appear that each signal-catcher was interrupted by the next
3883 signal. In this case, the implementation should ensure that this stacking of signals does not
3884 violate the semantics of the signal masks established by *sigaction()*. Other implementations
3885 process at most one signal when the operating system is entered, with remaining signals saved
3886 for later delivery. Although this practice is widespread, this behavior is neither standardized
3887 nor endorsed. In either case, implementations should attempt to deliver signals associated with
3888 the current state of the process (for example, SIGFPE) before other signals, if possible.

3889 In 4.2 BSD and 4.3 BSD, it is not permissible to ignore or explicitly block SIGCONT, because if
3890 blocking or ignoring this signal prevented it from continuing a stopped process, such a process
3891 could never be continued (only killed by SIGKILL). However, 4.2 BSD and 4.3 BSD do block
3892 SIGCONT during execution of its signal-catching function when it is caught, creating exactly
3893 this problem. A proposal was considered to disallow catching SIGCONT in addition to ignoring
3894 and blocking it, but this limitation led to objections. The consensus was to require that
3895 SIGCONT always continue a stopped process when generated. This removed the need to

3896 disallow ignoring or explicit blocking of the signal; note that SIG_IGN and SIG_DFL are
3897 equivalent for SIGCONT.

3898 B.2.4.2 Realtime Signal Generation and Delivery

3899 The realtime signals functionality is required in this version of the standard for the following
3900 reasons:

3901 • The **sigevent** structure is used by other POSIX.1 functions that result in asynchronous
3902 event notifications to specify the notification mechanism to use and other information
3903 needed by the notification mechanism. IEEE Std 1003.1-200x defines only three symbolic
3904 values for the notification mechanism:

3905 — SIGEV_NONE is used to indicate that no notification is required when the event
3906 occurs. This is useful for applications that use asynchronous I/O with polling for
3907 completion.

3908 — SIGEV_SIGNAL indicates that a signal is generated when the event occurs.

3909 — SIGEV_THREAD provides for “callback functions” for asynchronous notifications
3910 done by a function call within the context of a new thread. This provides a multi-
3911 threaded process with a more natural means of notification than signals.

3912 The primary difficulty with previous notification approaches has been to specify the
3913 environment of the notification routine.

3914 — One approach is to limit the notification routine to call only functions permitted in a
3915 signal handler. While the list of permissible functions is clearly stated, this is overly
3916 restrictive.

3917 — A second approach is to define a new list of functions or classes of functions that are
3918 explicitly permitted or not permitted. This would give a programmer more lists to
3919 deal with, which would be awkward.

3920 — The third approach is to define completely the environment for execution of the
3921 notification function. A clear definition of an execution environment for notification
3922 is provided by executing the notification function in the environment of a newly
3923 created thread.

3924 Implementations may support additional notification mechanisms by defining new values
3925 for *sigev_notify*.

3926 For a notification type of SIGEV_SIGNAL, the other members of the **sigevent** structure
3927 defined by IEEE Std 1003.1-200x specify the realtime signal—that is, the signal number and
3928 application-defined value that differentiates between occurrences of signals with the same
3929 number—that will be generated when the event occurs. The structure is defined in
3930 **<signal.h>**, even though the structure is not directly used by any of the signal functions,
3931 because it is part of the signals interface used by the POSIX.1b “client functions”. When the
3932 client functions include **<signal.h>** to define the signal names, the **sigevent** structure will
3933 also be defined.

3934 An application-defined value passed to the signal handler is used to differentiate between
3935 different “events” instead of requiring that the application use different signal numbers for
3936 several reasons:

3937 — Realtime applications potentially handle a very large number of different events.
3938 Requiring that implementations support a correspondingly large number of distinct
3939 signal numbers will adversely impact the performance of signal delivery because the
3940 signal masks to be manipulated on entry and exit to the handlers will become large.

- 3941 — Event notifications are prioritized by signal number (the rationale for this is
3942 explained in the following paragraphs) and the use of different signal numbers to
3943 differentiate between the different event notifications overloads the signal number
3944 more than has already been done. It also requires that the application writer make
3945 arbitrary assignments of priority to events that are logically of equal priority.

3946 A union is defined for the application-defined value so that either an integer constant or a
3947 pointer can be portably passed to the signal-catching function. On some architectures a
3948 pointer cannot be cast to an **int** and *vice versa*.

3949 Use of a structure here with an explicit notification type discriminant rather than explicit
3950 parameters to realtime functions, or embedded in other realtime structures, provides for
3951 future extensions to IEEE Std 1003.1-200x. Additional, perhaps more efficient, notification
3952 mechanisms can be supported for existing realtime function interfaces, such as timers and
3953 asynchronous I/O, by extending the **sigevent** structure appropriately. The existing
3954 realtime function interfaces will not have to be modified to use any such new notification
3955 mechanism. The revised text concerning the SIGEV_SIGNAL value makes consistent the
3956 semantics of the members of the **sigevent** structure, particularly in the definitions of
3957 *lio_listio()* and *aio_fsync()*. For uniformity, other revisions cause this specification to be
3958 referred to rather than inaccurately duplicated in the descriptions of functions and
3959 structures using the **sigevent** structure. The revised wording does not relax the
3960 requirement that the signal number be in the range SIGRTMIN to SIGRTMAX to guarantee
3961 queuing and passing of the application value, since that requirement is still implied by the
3962 signal names.

- 3963 • IEEE Std 1003.1-200x is intentionally vague on whether “non-realtime” signal-generating
3964 mechanisms can result in a **siginfo_t** being supplied to the handler on delivery. In one
3965 existing implementation, a **siginfo_t** is posted on signal generation, even though the
3966 implementation does not support queuing of multiple occurrences of a signal. It is not the
3967 intent of IEEE Std 1003.1-200x to preclude this, independent of the mandate to define
3968 signals that do support queuing. Any interpretation that appears to preclude this is a
3969 mistake in the reading or writing of the standard.
- 3970 • Signals handled by realtime signal handlers might be generated by functions or conditions
3971 that do not allow the specification of an application-defined value and do not queue.
3972 IEEE Std 1003.1-200x specifies the *si_code* member of the **siginfo_t** structure used in
3973 existing practice and defines additional codes so that applications can detect whether an
3974 application-defined value is present or not. The code SI_USER for *kill()*-generated signals
3975 is adopted from existing practice.
- 3976 • The *sigaction()* *sa_flags* value SA_SIGINFO tells the implementation that the signal-
3977 catching function expects two additional arguments. When the flag is not set, a single
3978 argument, the signal number, is passed as specified by IEEE Std 1003.1-200x. Although
3979 IEEE Std 1003.1-200x does not explicitly allow the *info* argument to the handler function to
3980 be NULL, this is existing practice. This provides for compatibility with programs whose
3981 signal-catching functions are not prepared to accept the additional arguments.
3982 IEEE Std 1003.1-200x is explicitly unspecified as to whether signals actually queue when
3983 SA_SIGINFO is not set for a signal, as there appear to be no benefits to applications in
3984 specifying one behavior or another. One existing implementation queues a **siginfo_t** on
3985 each signal generation, unless the signal is already pending, in which case the
3986 implementation discards the new **siginfo_t**; that is, the queue length is never greater than
3987 one. This implementation only examines SA_SIGINFO on signal delivery, discarding the
3988 queued **siginfo_t** if its delivery was not requested.

3989 IEEE Std 1003.1-200x specifies several new values for the *si_code* member of the **siginfo_t**
3990 structure. In existing practice, a *si_code* value of less than or equal to zero indicates that the

3991 signal was generated by a process via the `kill()` function. In existing practice, values of
 3992 `si_code` that provide additional information for implementation-generated signals, such as
 3993 SIGFPE or SIGSEGV, are all positive. Thus, if implementations define the new constants
 3994 specified in IEEE Std 1003.1-200x to be negative numbers, programs written to use existing
 3995 practice will not break. IEEE Std 1003.1-200x chose not to attempt to specify existing
 3996 practice values of `si_code` other than SI_USER both because it was deemed beyond the
 3997 scope of IEEE Std 1003.1-200x and because many of the values in existing practice appear
 3998 to be platform and implementation-defined. But, IEEE Std 1003.1-200x does specify that if
 3999 an implementation—for example, one that does not have existing practice in this area—
 4000 chooses to define additional values for `si_code`, these values have to be different from the
 4001 values of the symbols specified by IEEE Std 1003.1-200x. This will allow conforming
 4002 applications to differentiate between signals generated by one of the POSIX.1b
 4003 asynchronous events and those generated by other implementation events in a manner
 4004 compatible with existing practice.

4005 The unique values of `si_code` for the POSIX.1b asynchronous events have implications for
 4006 implementations of, for example, asynchronous I/O or message passing in user space
 4007 library code. Such an implementation will be required to provide a hidden interface to the
 4008 signal generation mechanism that allows the library to specify the standard values of
 4009 `si_code`.

4010 Existing practice also defines additional members of `siginfo_t`, such as the process ID and
 4011 user ID of the sending process for `kill()`-generated signals. These members were deemed
 4012 not necessary to meet the requirements of realtime applications and are not specified by
 4013 IEEE Std 1003.1-200x. Neither are they precluded.

4014 The third argument to the signal-catching function, `context`, is left undefined by
 4015 IEEE Std 1003.1-200x, but is specified in the interface because it matches existing practice
 4016 for the SA_SIGINFO flag. It was considered undesirable to require a separate
 4017 implementation for SA_SIGINFO for POSIX conformance on implementations that already
 4018 support the two additional parameters.

- 4019 • The requirement to deliver lower numbered signals in the range SIGRTMIN to SIGRTMAX
 4020 first, when multiple unblocked signals are pending, results from several considerations:

- 4021 — A method is required to prioritize event notifications. The signal number was chosen
 4022 instead of, for instance, associating a separate priority with each request, because an
 4023 implementation has to check pending signals at various points and select one for
 4024 delivery when more than one is pending. Specifying a selection order is the minimal
 4025 additional semantic that will achieve prioritized delivery. If a separate priority were
 4026 to be associated with queued signals, it would be necessary for an implementation to
 4027 search all non-empty, non-blocked signal queues and select from among them the
 4028 pending signal with the highest priority. This would significantly increase the cost of
 4029 and decrease the determinism of signal delivery.

- 4030 — Given the specified selection of the lowest numeric unblocked pending signal,
 4031 preemptive priority signal delivery can be achieved using signal numbers and signal
 4032 masks by ensuring that the `sa_mask` for each signal number blocks all signals with a
 4033 higher numeric value.

4034 For realtime applications that want to use only the newly defined realtime signal
 4035 numbers without interference from the standard signals, this can be achieved by
 4036 blocking all of the standard signals in the thread signal mask and in the `sa_mask`
 4037 installed by the signal action for the realtime signal handlers.

4038 IEEE Std 1003.1-200x explicitly leaves unspecified the ordering of signals outside of the
 4039 range of realtime signals and the ordering of signals within this range with respect to those

4040 outside the range. It was believed that this would unduly constrain implementations or
4041 standards in the future definition of new signals.

4042 B.2.4.3 Signal Actions

4043 Early proposals mentioned SIGCONT as a second exception to the rule that signals are not
4044 delivered to stopped processes until continued. Because IEEE Std 1003.1-200x now specifies that
4045 SIGCONT causes the stopped process to continue when it is generated, delivery of SIGCONT is
4046 not prevented because a process is stopped, even without an explicit exception to this rule.

4047 Ignoring a signal by setting the action to SIG_IGN (or SIG_DFL for signals whose default action
4048 is to ignore) is not the same as installing a signal-catching function that simply returns. Invoking
4049 such a function will interrupt certain system functions that block processes (for example, *wait()*,
4050 *sigsuspend()*, *pause()*, *read()*, *write()*) while ignoring a signal has no such effect on the process.

4051 Historical implementations discard pending signals when the action is set to SIG_IGN.
4052 However, they do not always do the same when the action is set to SIG_DFL and the default
4053 action is to ignore the signal. IEEE Std 1003.1-200x requires this for the sake of consistency and
4054 also for completeness, since the only signal this applies to is SIGCHLD, and
4055 IEEE Std 1003.1-200x disallows setting its action to SIG_IGN.

4056 Some implementations (System V, for example) assign different semantics for SIGCLD
4057 depending on whether the action is set to SIG_IGN or SIG_DFL. Since POSIX.1 requires that the
4058 default action for SIGCHLD be to ignore the signal, applications should always set the action to
4059 SIG_DFL in order to avoid SIGCHLD.

4060 Whether or not an implementation allows SIG_IGN as a SIGCHLD disposition to be inherited
4061 across a call to one of the *exec* family of functions or *posix_spawn()* is explicitly left as
4062 unspecified. This change was made as a result of IEEE PASC Interpretation 1003.1 #132, and
4063 permits the implementation to decide between the following alternatives:

- 4064 • Unconditionally leave SIGCHLD set to SIG_IGN, in which case the implementation would
4065 not allow applications that assume inheritance of SIG_DFL to conform to
4066 IEEE Std 1003.1-200x without change. The implementation would, however, retain an
4067 ability to control applications that create child processes but never call on the *wait* family of
4068 functions, potentially filling up the process table.
- 4069 • Unconditionally reset SIGCHLD to SIG_DFL, in which case the implementation would
4070 allow applications that assume inheritance of SIG_DFL to conform. The implementation
4071 would, however, lose an ability to control applications that spawn child processes but
4072 never reap them.
- 4073 • Provide some mechanism, not specified in IEEE Std 1003.1-200x, to control inherited
4074 SIGCHLD dispositions.

4075 Some implementations (System V, for example) will deliver a SIGCLD signal immediately when
4076 a process establishes a signal-catching function for SIGCLD when that process has a child that
4077 has already terminated. Other implementations, such as 4.3 BSD, do not generate a new
4078 SIGCHLD signal in this way. In general, a process should not attempt to alter the signal action
4079 for the SIGCHLD signal while it has any outstanding children. However, it is not always
4080 possible for a process to avoid this; for example, shells sometimes start up processes in pipelines
4081 with other processes from the pipeline as children. Processes that cannot ensure that they have
4082 no children when altering the signal action for SIGCHLD thus need to be prepared for, but not
4083 depend on, generation of an immediate SIGCHLD signal.

4084 The default action of the stop signals (SIGSTOP, SIGTSTP, SIGTTIN, SIGTTOU) is to stop a
4085 process that is executing. If a stop signal is delivered to a process that is already stopped, it has
4086 no effect. In fact, if a stop signal is generated for a stopped process whose signal mask blocks the

4087 signal, the signal will never be delivered to the process since the process must receive a
4088 SIGCONT, which discards all pending stop signals, in order to continue executing.

4089 The SIGCONT signal continues a stopped process even if SIGCONT is blocked (or ignored).
4090 However, if a signal-catching routine has been established for SIGCONT, it will not be entered
4091 until SIGCONT is unblocked.

4092 If a process in an orphaned process group stops, it is no longer under the control of a job control
4093 shell and hence would not normally ever be continued. Because of this, orphaned processes that
4094 receive terminal-related stop signals (SIGTSTP, SIGTTIN, SIGTTOU, but not SIGSTOP) must not
4095 be allowed to stop. The goal is to prevent stopped processes from languishing forever. (As
4096 SIGSTOP is sent only via *kill()*, it is assumed that the process or user sending a SIGSTOP can
4097 send a SIGCONT when desired.) Instead, the system must discard the stop signal. As an
4098 extension, it may also deliver another signal in its place. 4.3 BSD sends a SIGKILL, which is
4099 overly effective because SIGKILL is not catchable. Another possible choice is SIGHUP. 4.3 BSD
4100 also does this for orphaned processes (processes whose parent has terminated) rather than for
4101 members of orphaned process groups; this is less desirable because job control shells manage
4102 process groups. POSIX.1 also prevents SIGTTIN and SIGTTOU signals from being generated for
4103 processes in orphaned process groups as a direct result of activity on a terminal, preventing
4104 infinite loops when *read()* and *write()* calls generate signals that are discarded; see [Section](#)
4105 [A.11.1.4](#) (on page 68). A similar restriction on the generation of SIGTSTP was considered, but
4106 that would be unnecessary and more difficult to implement due to its asynchronous nature.

4107 Although POSIX.1 requires that signal-catching functions be called with only one argument,
4108 there is nothing to prevent conforming implementations from extending POSIX.1 to pass
4109 additional arguments, as long as Strictly Conforming POSIX.1 Applications continue to compile
4110 and execute correctly. Most historical implementations do, in fact, pass additional, signal-
4111 specific arguments to certain signal-catching routines.

4112 There was a proposal to change the declared type of the signal handler to:

```
4113 void func (int sig, ...);
```

4114 The usage of ellipses ("...") is ISO C standard syntax to indicate a variable number of
4115 arguments. Its use was intended to allow the implementation to pass additional information to
4116 the signal handler in a standard manner.

4117 Unfortunately, this construct would require all signal handlers to be defined with this syntax
4118 because the ISO C standard allows implementations to use a different parameter passing
4119 mechanism for variable parameter lists than for non-variable parameter lists. Thus, all existing
4120 signal handlers in all existing applications would have to be changed to use the variable syntax
4121 in order to be standard and portable. This is in conflict with the goal of Minimal Changes to
4122 Existing Application Code.

4123 When terminating a process from a signal-catching function, processes should be aware of any
4124 interpretation that their parent may make of the status returned by *wait()* or *waitpid()*. In
4125 particular, a signal-catching function should not call *exit(0)* or *_exit(0)* unless it wants to indicate
4126 successful termination. A non-zero argument to *exit()* or *_exit()* can be used to indicate
4127 unsuccessful termination. Alternatively, the process can use *kill()* to send itself a fatal signal
4128 (first ensuring that the signal is set to the default action and not blocked). See also the
4129 RATIONALE section of the *_exit()* function.

4130 The behavior of *unsafe* functions, as defined by this section, is undefined when they are invoked
4131 from signal-catching functions in certain circumstances. The behavior of reentrant functions, as
4132 defined by this section, is as specified by POSIX.1, regardless of invocation from a signal-
4133 catching function. This is the only intended meaning of the statement that reentrant functions
4134 may be used in signal-catching functions without restriction. Applications must still consider all

4135 effects of such functions on such things as data structures, files, and process state. In particular,
 4136 application writers need to consider the restrictions on interactions when interrupting *sleep()*
 4137 (see *sleep()*) and interactions among multiple handles for a file description. The fact that any
 4138 specific function is listed as reentrant does not necessarily mean that invocation of that function
 4139 from a signal-catching function is recommended.

4140 In order to prevent errors arising from interrupting non-reentrant function calls, applications
 4141 should protect calls to these functions either by blocking the appropriate signals or through the
 4142 use of some programmatic semaphore. POSIX.1 does not address the more general problem of
 4143 synchronizing access to shared data structures. Note in particular that even the “safe” functions
 4144 may modify the global variable *errno*; the signal-catching function may want to save and restore
 4145 its value. The same principles apply to the reentrancy of application routines and asynchronous
 4146 data access.

4147 Note that *longjmp()* and *siglongjmp()* are not in the list of reentrant functions. This is because the
 4148 code executing after *longjmp()* or *siglongjmp()* can call any unsafe functions with the same
 4149 danger as calling those unsafe functions directly from the signal handler. Applications that use
 4150 *longjmp()* or *siglongjmp()* out of signal handlers require rigorous protection in order to be
 4151 portable. Many of the other functions that are excluded from the list are traditionally
 4152 implemented using either the C language *malloc()* or *free()* functions or the ISO C standard I/O
 4153 library, both of which traditionally use data structures in a non-reentrant manner. Because any
 4154 combination of different functions using a common data structure can cause reentrancy
 4155 problems, POSIX.1 does not define the behavior when any unsafe function is called in a signal
 4156 handler that interrupts any unsafe function.

4157 The only realtime extension to signal actions is the addition of the additional parameters to the
 4158 signal-catching function. This extension has been explained and motivated in the previous
 4159 section. In making this extension, though, developers of POSIX.1b ran into issues relating to
 4160 function prototypes. In response to input from the POSIX.1 standard developers, members were
 4161 added to the **sigaction** structure to specify function prototypes for the newer signal-catching
 4162 function specified by POSIX.1b. These members follow changes that are being made to POSIX.1.
 4163 Note that IEEE Std 1003.1-200x explicitly states that these fields may overlap so that a union can
 4164 be defined. This enabled existing implementations of POSIX.1 to maintain binary-compatibility
 4165 when these extensions were added.

4166 The **siginfo_t** structure was adopted for passing the application-defined value to match existing
 4167 practice, but the existing practice has no provision for an application-defined value, so this was
 4168 added. Note that POSIX normally reserves the “_t” type designation for opaque types. The
 4169 **siginfo_t** structure breaks with this convention to follow existing practice and thus promote
 4170 portability. Standardization of the existing practice for the other members of this structure may
 4171 be addressed in the future.

4172 Although it is not explicitly visible to applications, there are additional semantics for signal
 4173 actions implied by queued signals and their interaction with other POSIX.1b realtime functions.
 4174 Specifically:

- 4175 • It is not necessary to queue signals whose action is SIG_IGN.
- 4176 • For implementations that support POSIX.1b timers, some interaction with the timer
 4177 functions at signal delivery is implied to manage the timer overrun count.

4178 IEEE Std 1003.1-2001/Cor 1-2002, item XSH/TC1/D6/5 is applied, reordering the RTS shaded
 4179 text under the third and fourth paragraphs of the SIG_DFL description. This corrects an earlier
 4180 editorial error in this section.

4181 IEEE Std 1003.1-2001/Cor 1-2002, item XSH/TC1/D6/6 is applied, adding the *abort()* function
 4182 to the list of async-cancel-safe functions.

4183 IEEE Std 1003.1-2001/Cor 2-2004, item XSH/TC2/D6/4 is applied, adding the *socketmark()*
 4184 function to the list of functions that shall be either reentrant or non-interruptible by signals and
 4185 shall be async-signal-safe.

4186 B.2.4.4 *Signal Effects on Other Functions*

4187 The most common behavior of an interrupted function after a signal-catching function returns is
 4188 for the interrupted function to give an [EINTR] error unless the SA_RESTART flag is in effect for
 4189 the signal. However, there are a number of specific exceptions, including *sleep()* and certain
 4190 situations with *read()* and *write()*.

4191 The historical implementations of many functions defined by IEEE Std 1003.1-200x are not
 4192 interruptible, but delay delivery of signals generated during their execution until after they
 4193 complete. This is never a problem for functions that are guaranteed to complete in a short
 4194 (imperceptible to a human) period of time. It is normally those functions that can suspend a
 4195 process indefinitely or for long periods of time (for example, *wait()*, *pause()*, *sigsuspend()*,
 4196 *sleep()*, or *read()/write()* on a slow device like a terminal) that are interruptible. This permits
 4197 applications to respond to interactive signals or to set timeouts on calls to most such functions
 4198 with *alarm()*. Therefore, implementations should generally make such functions (including ones
 4199 defined as extensions) interruptible.

4200 Functions not mentioned explicitly as interruptible may be so on some implementations,
 4201 possibly as an extension where the function gives an [EINTR] error. There are several functions
 4202 (for example, *getpid()*, *getuid()*) that are specified as never returning an error, which can thus
 4203 never be extended in this way.

4204 If a signal-catching function returns while the SA_RESTART flag is in effect, an interrupted
 4205 function is restarted at the point it was interrupted. Conforming applications cannot make
 4206 assumptions about the internal behavior of interrupted functions, even if the functions are
 4207 async-signal-safe. For example, suppose the *read()* function is interrupted with SA_RESTART in
 4208 effect, the signal-catching function closes the file descriptor being read from and returns, and the
 4209 *read()* function is then restarted; in this case the application cannot assume that the *read()*
 4210 function will give an [EBADF] error, since *read()* might have checked the file descriptor for
 4211 validity before being interrupted.

4212 B.2.5 Standard I/O Streams

4213 B.2.5.1 *Interaction of File Descriptors and Standard I/O Streams*

4214 There is no additional rationale provided for this section.

4215 B.2.5.2 *Stream Orientation and Encoding Rules*

4216 There is no additional rationale provided for this section.

4217 B.2.6 STREAMS

4218 STREAMS are included into IEEE Std 1003.1-200x as part of the alignment with the Single UNIX
 4219 Specification, but marked as an option in recognition that not all systems may wish to
 4220 implement the facility. The option within IEEE Std 1003.1-200x is denoted by the XSR margin
 4221 marker. The standard developers made this option independent of the XSI option. In this
 4222 version of the standard this option is marked obsolescent.

4223 STREAMS are a method of implementing network services and other character-based
 4224 input/output mechanisms, with the STREAM being a full-duplex connection between a process
 4225 and a device. STREAMS provides direct access to protocol modules, and optional protocol
 4226 modules can be interposed between the process-end of the STREAM and the device-driver at the

4227 device-end of the STREAM. Pipes can be implemented using the STREAMS mechanism, so they
4228 can provide process-to-process as well as process-to-device communications.

4229 This section introduces STREAMS I/O, the message types used to control them, an overview of
4230 the priority mechanism, and the interfaces used to access them.

4231 B.2.6.1 Accessing STREAMS

4232 There is no additional rationale provided for this section.

4233 B.2.7 XSI Interprocess Communication

4234 There are two forms of IPC supported as options in IEEE Std 1003.1-200x. The traditional
4235 System V IPC routines derived from the SVID—that is, the *msg**(), *sem**(), and *shm**()
4236 interfaces—are mandatory on XSI-conformant systems. Thus, all XSI-conformant systems
4237 provide the same mechanisms for manipulating messages, shared memory, and semaphores.

4238 In addition, the POSIX Realtime Extension provides an alternate set of routines for those systems
4239 supporting the appropriate options.

4240 The application writer is presented with a choice: the System V interfaces or the POSIX
4241 interfaces (loosely derived from the Berkeley interfaces). The XSI profile prefers the System V
4242 interfaces, but the POSIX interfaces may be more suitable for realtime or other performance-
4243 sensitive applications.

4244 B.2.7.1 IPC General Information

4245 General information that is shared by all three mechanisms is described in this section. The
4246 common permissions mechanism is briefly introduced, describing the mode bits, and how they
4247 are used to determine whether or not a process has access to read or write/alter the appropriate
4248 instance of one of the IPC mechanisms. All other relevant information is contained in the
4249 reference pages themselves.

4250 The semaphore type of IPC allows processes to communicate through the exchange of
4251 semaphore values. A semaphore is a positive integer. Since many applications require the use of
4252 more than one semaphore, XSI-conformant systems have the ability to create sets or arrays of
4253 semaphores.

4254 Calls to support semaphores include:

4255 *semctl*(), *semget*(), *semop*()

4256 Semaphore sets are created by using the *semget*() function.

4257 The message type of IPC allows processes to communicate through the exchange of data stored
4258 in buffers. This data is transmitted between processes in discrete portions known as messages.

4259 Calls to support message queues include:

4260 *msgctl*(), *msgget*(), *msgrcv*(), *msgsnd*()

4261 The shared memory type of IPC allows two or more processes to share memory and
4262 consequently the data contained therein. This is done by allowing processes to set up access to a
4263 common memory address space. This sharing of memory provides a fast means of exchange of
4264 data between processes.

4265 Calls to support shared memory include:

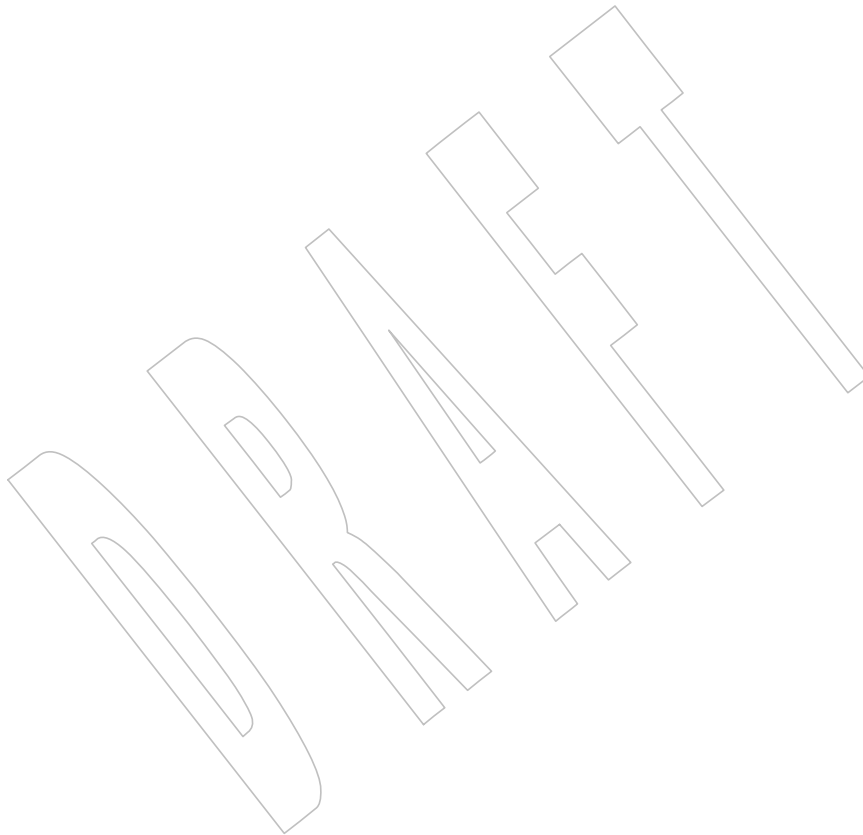
4266 *shmctl*(), *shmdt*(), *shmget*()

4267 The *ftok*() interface is also provided.

B.2.8 Realtime

Advisory Information

POSIX.1b contains an Informative Annex with proposed interfaces for “realtime files”. These interfaces could determine groups of the exact parameters required to do “direct I/O” or “extents”. These interfaces were objected to by a significant portion of the balloting group as too complex. A conforming application had little chance of correctly navigating the large parameter space to match its desires to the system. In addition, they only applied to a new type of file (realtime files) and they told the implementation exactly what to do as opposed to advising the implementation on application behavior and letting it optimize for the system the (portable)



4313 Implementations may use information conveyed by a previous *posix_fadvise()* call to influence
 4314 the manner in which allocation is performed. For example, if an application did the following
 4315 calls:

```
4316 fd = open("file");
4317 posix_fadvise(fd, offset, len, POSIX_FADV_SEQUENTIAL);
4318 posix_fallocate(fd, len, size);
```

4319 an implementation might allocate the file contiguously on disk.

4320 Finally, the *pathconf()* variables {POSIX_REC_MIN_XFER_SIZE},
 4321 {POSIX_REC_MAX_XFER_SIZE}, and {POSIX_REC_INCR_XFER_SIZE} tell the application a
 4322 range of transfer sizes that are recommended for best I/O performance.

4323 Where bounded response time is required, the vendor can supply the appropriate settings of the
 4324 advisories to achieve a guaranteed performance level.

4325 The interfaces meet the goals while allowing applications using regular files to take advantage
 4326 of performance optimizations. The interfaces tell the implementation expected application
 4327 behavior which the implementation can use to optimize performance on a particular system
 4328 with a particular dynamic load.

4329 The *posix_memalign()* function was added to allow for the allocation of specifically aligned
 4330 buffers; for example, for {POSIX_REC_XFER_ALIGN}.

4331 The working group also considered the alternative of adding a function which would return an
 4332 aligned pointer to memory within a user-supplied buffer. This was not considered to be the best
 4333 method, because it potentially wastes large amounts of memory when buffers need to be aligned
 4334 on large alignment boundaries.

4335 Message Passing

4336 This section provides the rationale for the definition of the message passing interface in
 4337 IEEE Std 1003.1-200x. This is presented in terms of the objectives, models, and requirements
 4338 imposed upon this interface.

- 4339 • Objectives

4340 Many applications, including both realtime and database applications, require a means of
 4341 passing arbitrary amounts of data between cooperating processes comprising the overall
 4342 application on one or more processors. Many conventional interfaces for interprocess
 4343 communication are insufficient for realtime applications in that efficient and deterministic
 4344 data passing methods cannot be implemented. This has prompted the definition of
 4345 message passing interfaces providing these facilities:

- 4346 — Open a message queue.
- 4347 — Send a message to a message queue.
- 4348 — Receive a message from a queue, either synchronously or asynchronously.
- 4349 — Alter message queue attributes for flow and resource control.

4350 It is assumed that an application may consist of multiple cooperating processes and that
 4351 these processes may wish to communicate and coordinate their activities. The message
 4352 passing facility described in IEEE Std 1003.1-200x allows processes to communicate
 4353 through system-wide queues. These message queues are accessed through names that
 4354 may be pathnames. A message queue can be opened for use by multiple sending and/or
 4355 multiple receiving processes.

4356 • Background on Embedded Applications

4357 Interprocess communication utilizing message passing is a key facility for the construction
 4358 of deterministic, high-performance realtime applications. The facility is present in all
 4359 realtime systems and is the framework upon which the application is constructed. The
 4360 performance of the facility is usually a direct indication of the performance of the resulting
 4361 application.

4362 Realtime applications, especially for embedded systems, are typically designed around the
 4363 performance constraints imposed by the message passing mechanisms. Applications for
 4364 embedded systems are typically very tightly constrained. Application writers expect to
 4365 design and control the entire system. In order to minimize system costs, the writer will
 4366 attempt to use all resources to their utmost and minimize the requirement to add
 4367 additional memory or processors.

4368 The embedded applications usually share address spaces and only a simple message
 4369 passing mechanism is required. The application can readily access common data incurring
 4370 only mutual-exclusion overheads. The models desired are the simplest possible with the
 4371 application building higher-level facilities only when needed.

4372 • Requirements

4373 The following requirements determined the features of the message passing facilities
 4374 defined in IEEE Std 1003.1-200x:

4375 — Naming of Message Queues

4376 The mechanism for gaining access to a message queue is a pathname evaluated in a
 4377 context that is allowed to be a file system name space, or it can be independent of
 4378 any file system. This is a specific attempt to allow implementations based on either
 4379 method in order to address both embedded systems and to also allow
 4380 implementation in larger systems.

4381 The interface of `mq_open()` is defined to allow but not require the access control and
 4382 name conflicts resulting from utilizing a file system for name resolution. All required
 4383 behavior is specified for the access control case. Yet a conforming implementation,
 4384 such as an embedded system kernel, may define that there are no distinctions
 4385 between users and may define that all processes have all access privileges.

4386 — Embedded System Naming

4387 Embedded systems need to be able to utilize independent name spaces for accessing
 4388 the various system objects. They typically do not have a file system, precluding its
 4389 utilization as a common name resolution mechanism. The modularity of an
 4390 embedded system limits the connections between separate mechanisms that can be
 4391 allowed.

4392 Embedded systems typically do not have any access protection. Since the system
 4393 does not support the mixing of applications from different areas, and usually does
 4394 not even have the concept of an authorization entity, access control is not useful.

4395 — Large System Naming

4396 On systems with more functionality, the name resolution must support the ability to
 4397 use the file system as the name resolution mechanism/object storage medium and to
 4398 have control over access to the objects. Utilizing the pathname space can result in
 4399 further errors when the names conflict with other objects.

- 4400 — Fixed Size of Messages
- 4401 The interfaces impose a fixed upper bound on the size of messages that can be sent to
4402 a specific message queue. The size is set on an individual queue basis and cannot be
4403 changed dynamically.
- 4404 The purpose of the fixed size is to increase the ability of the system to optimize the
4405 implementation of *mq_send()* and *mq_receive()*. With fixed sizes of messages and
4406 fixed numbers of messages, specific message blocks can be pre-allocated. This
4407 eliminates a significant amount of checking for errors and boundary conditions.
4408 Additionally, an implementation can optimize data copying to maximize
4409 performance. Finally, with a restricted range of message sizes, an implementation is
4410 better able to provide deterministic operations.
- 4411 — Prioritization of Messages
- 4412 Message prioritization allows the application to determine the order in which
4413 messages are received. Prioritization of messages is a key facility that is provided by
4414 most realtime kernels and is heavily utilized by the applications. The major purpose
4415 of having priorities in message queues is to avoid priority inversions in the message
4416 system, where a high-priority message is delayed behind one or more lower-priority
4417 messages. This allows the applications to be designed so that they do not need to be
4418 interrupted in order to change the flow of control when exceptional conditions occur.
4419 The prioritization does add additional overhead to the message operations in those
4420 cases it is actually used but a clever implementation can optimize for the FIFO case to
4421 make that more efficient.
- 4422 — Asynchronous Notification
- 4423 The interface supports the ability to have a task asynchronously notified of the
4424 availability of a message on the queue. The purpose of this facility is to allow the task
4425 to perform other functions and yet still be notified that a message has become
4426 available on the queue.
- 4427 To understand the requirement for this function, it is useful to understand two
4428 models of application design: a single task performing multiple functions and
4429 multiple tasks performing a single function. Each of these models has advantages.
- 4430 Asynchronous notification is required to build the model of a single task performing
4431 multiple operations. This model typically results from either the expectation that
4432 interruption is less expensive than utilizing a separate task or from the growth of the
4433 application to include additional functions.

4434 Semaphores

4435 Semaphores are a high-performance process synchronization mechanism. Semaphores are
4436 named by null-terminated strings of characters.

4437 A semaphore is created using the *sem_init()* function or the *sem_open()* function with the
4438 *O_CREAT* flag set in *oflag*.

4439 To use a semaphore, a process has to first initialize the semaphore or inherit an open descriptor
4440 for the semaphore via *fork()*.

4441 A semaphore preserves its state when the last reference is closed. For example, if a semaphore
4442 has a value of 13 when the last reference is closed, it will have a value of 13 when it is next
4443 opened.

4444 When a semaphore is created, an initial state for the semaphore has to be provided. This value is

4445 a non-negative integer. Negative values are not possible since they indicate the presence of
 4446 blocked processes. The persistence of any of these objects across a system crash or a system
 4447 reboot is undefined. Conforming applications must not depend on any sort of persistence across
 4448 a system reboot or a system crash.

4449 • Models and Requirements

4450 A realtime system requires synchronization and communication between the processes
 4451 comprising the overall application. An efficient and reliable synchronization mechanism
 4452 has to be provided in a realtime system that will allow more than one schedulable process
 4453 mutually-exclusive access to the same resource. This synchronization mechanism has to
 4454 allow for the optimal implementation of synchronization or systems implementors will
 4455 define other, more cost-effective methods.

4456 At issue are the methods whereby multiple processes (tasks) can be designed and
 4457 implemented to work together in order to perform a single function. This requires
 4458 interprocess communication and synchronization. A semaphore mechanism is the lowest
 4459 level of synchronization that can be provided by an operating system.

4460 A semaphore is defined as an object that has an integral value and a set of blocked
 4461 processes associated with it. If the value is positive or zero, then the set of blocked
 4462 processes is empty; otherwise, the size of the set is equal to the absolute value of the
 4463 semaphore value. The value of the semaphore can be incremented or decremented by any
 4464 process with access to the semaphore and must be done as an indivisible operation. When
 4465 a semaphore value is less than or equal to zero, any process that attempts to lock it again
 4466 will block or be informed that it is not possible to perform the operation.

4467 A semaphore may be used to guard access to any resource accessible by more than one
 4468 schedulable task in the system. It is a global entity and not associated with any particular
 4469 process. As such, a method of obtaining access to the semaphore has to be provided by the
 4470 operating system. A process that wants access to a critical resource (section) has to wait on
 4471 the semaphore that guards that resource. When the semaphore is locked on behalf of a
 4472 process, it knows that it can utilize the resource without interference by any other
 4473 cooperating process in the system. When the process finishes its operation on the resource,
 4474 leaving it in a well-defined state, it posts the semaphore, indicating that some other
 4475 process may now obtain the resource associated with that semaphore.

4476 In this section, mutexes and condition variables are specified as the synchronization
 4477 mechanisms between threads.

4478 These primitives are typically used for synchronizing threads that share memory in a
 4479 single process. However, this section provides an option allowing the use of these
 4480 synchronization interfaces and objects between processes that share memory, regardless of
 4481 the method for sharing memory.

4482 Much experience with semaphores shows that there are two distinct uses of
 4483 synchronization: locking, which is typically of short duration; and waiting, which is
 4484 typically of long or unbounded duration. These distinct usages map directly onto mutexes
 4485 and condition variables, respectively.

4486 Semaphores are provided in IEEE Std 1003.1-200x primarily to provide a means of
 4487 synchronization for processes; these processes may or may not share memory. Mutexes
 4488 and condition variables are specified as synchronization mechanisms between threads;
 4489 these threads always share (some) memory. Both are synchronization paradigms that have
 4490 been in widespread use for a number of years. Each set of primitives is particularly well
 4491 matched to certain problems.

4492 With respect to binary semaphores, experience has shown that condition variables and

4493 mutexes are easier to use for many synchronization problems than binary semaphores. The
 4494 primary reason for this is the explicit appearance of a Boolean predicate that specifies
 4495 when the condition wait is satisfied. This Boolean predicate terminates a loop, including
 4496 the call to *pthread_cond_wait()*. As a result, extra wakeups are benign since the predicate
 4497 governs whether the thread will actually proceed past the condition wait. With stateful
 4498 primitives, such as binary semaphores, the wakeup in itself typically means that the wait is
 4499 satisfied. The burden of ensuring correctness for such waits is thus placed on *all* signalers
 4500 of the semaphore rather than on an *explicitly coded* Boolean predicate located at the
 4501 condition wait. Experience has shown that the latter creates a major improvement in safety
 4502 and ease-of-use.

4503 Counting semaphores are well matched to dealing with producer/consumer problems,
 4504 including those that might exist between threads of different processes, or between a signal
 4505 handler and a thread. In the former case, there may be little or no memory shared by the
 4506 processes; in the latter case, one is not communicating between co-equal threads, but
 4507 between a thread and an interrupt-like entity. It is for these reasons that
 4508 IEEE Std 1003.1-200x allows semaphores to be used by threads.

4509 Mutexes and condition variables have been effectively used with and without priority
 4510 inheritance, priority ceiling, and other attributes to synchronize threads that share
 4511 memory. The efficiency of their implementation is comparable to or better than that of
 4512 other synchronization primitives that are sometimes harder to use (for example, binary
 4513 semaphores). Furthermore, there is at least one known implementation of Ada tasking that
 4514 uses these primitives. Mutexes and condition variables together constitute an appropriate,
 4515 sufficient, and complete set of inter-thread synchronization primitives.

4516 Efficient multi-threaded applications require high-performance synchronization
 4517 primitives. Considerations of efficiency and generality require a small set of primitives
 4518 upon which more sophisticated synchronization functions can be built.

4519 • Standardization Issues

4520 It is possible to implement very high-performance semaphores using test-and-set
 4521 instructions on shared memory locations. The library routines that implement such a high-
 4522 performance interface have to properly ensure that a *sem_wait()* or *sem_trywait()* operation
 4523 that cannot be performed will issue a blocking semaphore system call or properly report
 4524 the condition to the application. The same interface to the application program would be
 4525 provided by a high-performance implementation.

4526 **B.2.8.1 Realtime Signals**

4527 **Realtime Signals Extension**

4528 This portion of the rationale presents models, requirements, and standardization issues relevant
 4529 to the Realtime Signals Extension. This extension provides the capability required to support
 4530 reliable, deterministic, asynchronous notification of events. While a new mechanism,
 4531 unencumbered by the historical usage and semantics of POSIX.1 signals, might allow for a more
 4532 efficient implementation, the application requirements for event notification can be met with a
 4533 small number of extensions to signals. Therefore, a minimal set of extensions to signals to
 4534 support the application requirements is specified.

4535 The realtime signal extensions specified in this section are used by other realtime functions
 4536 requiring asynchronous notification:

4537 • Models

4538 The model supported is one of multiple cooperating processes, each of which handles
 4539 multiple asynchronous external events. Events represent occurrences that are generated as

4540 the result of some activity in the system. Examples of occurrences that can constitute an
4541 event include:

- 4542 — Completion of an asynchronous I/O request
- 4543 — Expiration of a POSIX.1b timer
- 4544 — Arrival of an interprocess message
- 4545 — Generation of a user-defined event

4546 Processing of these events may occur synchronously via polling for event notifications or
4547 asynchronously via a software interrupt mechanism. Existing practice for this model is
4548 well established for traditional proprietary realtime operating systems, realtime
4549 executives, and realtime extended POSIX-like systems.

4550 A contrasting model is that of “cooperating sequential processes” where each process
4551 handles a single priority of events via polling. Each process blocks while waiting for
4552 events, and each process depends on the preemptive, priority-based process scheduling
4553 mechanism to arbitrate between events of different priority that need to be processed
4554 concurrently. Existing practice for this model is also well established for small realtime
4555 executives that typically execute in an unprotected physical address space, but it is just
4556 emerging in the context of a fuller function operating system with multiple virtual address
4557 spaces.

4558 It could be argued that the cooperating sequential process model, and the facilities
4559 supported by the POSIX Threads Extension obviate a software interrupt model. But, even
4560 with the cooperating sequential process model, the need has been recognized for a
4561 software interrupt model to handle exceptional conditions and process aborting, so the
4562 mechanism must be supported in any case. Furthermore, it is not the purview of
4563 IEEE Std 1003.1-200x to attempt to convince realtime practitioners that their current
4564 application models based on software interrupts are “broken” and should be replaced by
4565 the cooperating sequential process model. Rather, it is the charter of IEEE Std 1003.1-200x
4566 to provide standard extensions to mechanisms that support existing realtime practice.

4567 • Requirements

4568 This section discusses the following realtime application requirements for asynchronous
4569 event notification:

- 4570 — Reliable delivery of asynchronous event notification

4571 The events notification mechanism guarantees delivery of an event notification.
4572 Asynchronous operations (such as asynchronous I/O and timers) that complete
4573 significantly after they are invoked have to guarantee that delivery of the event
4574 notification can occur at the time of completion.

- 4575 — Prioritized handling of asynchronous event notifications

4576 The events notification mechanism supports the assigning of a user function as an
4577 event notification handler. Furthermore, the mechanism supports the preemption of
4578 an event handler function by a higher priority event notification and supports the
4579 selection of the highest priority pending event notification when multiple
4580 notifications (of different priority) are pending simultaneously.

4581 The model here is based on hardware interrupts. Asynchronous event handling
4582 allows the application to ensure that time-critical events are immediately processed
4583 when delivered, without the indeterminism of being at a random location within a
4584 polling loop. Use of handler priority allows the specification of how handlers are
4585 interrupted by other higher priority handlers.

- 4586 — Differentiation between multiple occurrences of event notifications of the same type
- 4587 The events notification mechanism passes an application-defined value to the event
4588 handler function. This value can be used for a variety of purposes, such as enabling
4589 the application to identify which of several possible events of the same type (for
4590 example, timer expirations) has occurred.
- 4591 — Polled reception of asynchronous event notifications
- 4592 The events notification mechanism supports blocking and non-blocking polls for
4593 asynchronous event notification.
- 4594 The polled mode of operation is often preferred over the interrupt mode by those
4595 practitioners accustomed to this model. Providing support for this model facilitates
4596 the porting of applications based on this model to POSIX.1b conforming systems.
- 4597 — Deterministic response to asynchronous event notifications
- 4598 The events notification mechanism does not preclude implementations that provide
4599 deterministic event dispatch latency and minimizes the number of system calls
4600 needed to use the event facilities during realtime processing.
- 4601 • Rationale for Extension
- 4602 POSIX.1 signals have many of the characteristics necessary to support the asynchronous
4603 handling of event notifications, and the Realtime Signals Extension addresses the
4604 following deficiencies in the POSIX.1 signal mechanism:
- 4605 — Signals do not support reliable delivery of event notification. Subsequent
4606 occurrences of a pending signal are not guaranteed to be delivered.
 - 4607 — Signals do not support prioritized delivery of event notifications. The order of signal
4608 delivery when multiple unblocked signals are pending is undefined.
 - 4609 — Signals do not support the differentiation between multiple signals of the same type.

4610 B.2.8.2 Asynchronous I/O

4611 Many applications need to interact with the I/O subsystem in an asynchronous manner. The
4612 asynchronous I/O mechanism provides the ability to overlap application processing and I/O
4613 operations initiated by the application. The asynchronous I/O mechanism allows a single
4614 process to perform I/O simultaneously to a single file multiple times or to multiple files
4615 multiple times.

4616 Overview

4617 Asynchronous I/O operations proceed in logical parallel with the processing done by the
4618 application after the asynchronous I/O has been initiated. Other than this difference,
4619 asynchronous I/O behaves similarly to normal I/O using *read()*, *write()*, *lseek()*, and *fsync()*.
4620 The effect of issuing an asynchronous I/O request is as if a separate thread of execution were to
4621 perform atomically the implied *lseek()* operation, if any, and then the requested I/O operation
4622 (either *read()*, *write()*, or *fsync()*). There is no seek implied with a call to *aio_fsync()*. Concurrent
4623 asynchronous operations and synchronous operations applied to the same file update the file as
4624 if the I/O operations had proceeded serially.

4625 When asynchronous I/O completes, a signal can be delivered to the application to indicate the
4626 completion of the I/O. This signal can be used to indicate that buffers and control blocks used
4627 for asynchronous I/O can be reused. Signal delivery is not required for an asynchronous
4628 operation and may be turned off on a per-operation basis by the application. Signals may also be
4629 synchronously polled using *aio_suspend()*, *sigtimedwait()*, or *sigwaitinfo()*.

4630 Normal I/O has a return value and an error status associated with it. Asynchronous I/O
 4631 returns a value and an error status when the operation is first submitted, but that only relates to
 4632 whether the operation was successfully queued up for servicing. The I/O operation itself also
 4633 has a return status and an error value. To allow the application to retrieve the return status and
 4634 the error value, functions are provided that, given the address of an asynchronous I/O control
 4635 block, yield the return and error status associated with the operation. Until an asynchronous I/O
 4636 operation is done, its error status is [EINPROGRESS]. Thus, an application can poll for
 4637 completion of an asynchronous I/O operation by waiting for the error status to become equal to
 4638 a value other than [EINPROGRESS]. The return status of an asynchronous I/O operation is
 4639 undefined so long as the error status is equal to [EINPROGRESS].

4640 Storage for asynchronous operation return and error status may be limited. Submission of
 4641 asynchronous I/O operations may fail if this storage is exceeded. When an application retrieves
 4642 the return status of a given asynchronous operation, therefore, any system-maintained storage
 4643 used for this status and the error status may be reclaimed for use by other asynchronous
 4644 operations.

4645 Asynchronous I/O can be performed on file descriptors that have been enabled for POSIX.1b
 4646 synchronized I/O. In this case, the I/O operation still occurs asynchronously, as defined herein;
 4647 however, the asynchronous operation I/O in this case is not completed until the I/O has reached
 4648 either the state of synchronized I/O data integrity completion or synchronized I/O file integrity
 4649 completion, depending on the sort of synchronized I/O that is enabled on the file descriptor.

4650 Models

4651 Three models illustrate the use of asynchronous I/O: a journalization model, a data acquisition
 4652 model, and a model of the use of asynchronous I/O in supercomputing applications.

- 4653 • Journalization Model

4654 Many realtime applications perform low-priority journalizing functions. Journalizing
 4655 requires that logging records be queued for output without blocking the initiating process.

- 4656 • Data Acquisition Model

4657 A data acquisition process may also serve as a model. The process has two or more
 4658 channels delivering intermittent data that must be read within a certain time. The process
 4659 issues one asynchronous read on each channel. When one of the channels needs data
 4660 collection, the process reads the data and posts it through an asynchronous write to
 4661 secondary memory for future processing.

- 4662 • Supercomputing Model

4663 The supercomputing community has used asynchronous I/O much like that specified in
 4664 POSIX.1 for many years. This community requires the ability to perform multiple I/O
 4665 operations to multiple devices with a minimal number of entries to “the system”; each
 4666 entry to “the system” provokes a major delay in operations when compared to the normal
 4667 progress made by the application. This existing practice motivated the use of combined
 4668 *lseek()* and *read()* or *write()* calls, as well as the *lio_listio()* call. Another common practice is
 4669 to disable signal notification for I/O completion, and simply poll for I/O completion at
 4670 some interval by which the I/O should be completed. Likewise, interfaces like *aio_cancel()*
 4671 have been in successful commercial use for many years. Note also that an underlying
 4672 implementation of asynchronous I/O will require the ability, at least internally, to cancel
 4673 outstanding asynchronous I/O, at least when the process exits. (Consider an asynchronous
 4674 read from a terminal, when the process intends to exit immediately.)

4675 **Requirements**

4676 Asynchronous input and output for realtime implementations have these requirements:

- 4677 • The ability to queue multiple asynchronous read and write operations to a single open
- 4678 instance. Both sequential and random access should be supported.
- 4679 • The ability to queue asynchronous read and write operations to multiple open instances.
- 4680 • The ability to obtain completion status information by polling and/or asynchronous event
- 4681 notification.
- 4682 • Asynchronous event notification on asynchronous I/O completion is optional.
- 4683 • It has to be possible for the application to associate the event with the *aiocbp* for the
- 4684 operation that generated the event.
- 4685 • The ability to cancel queued requests.
- 4686 • The ability to wait upon asynchronous I/O completion in conjunction with other types of
- 4687 events.
- 4688 • The ability to accept an *aio_read()* and an *aio_cancel()* for a device that accepts a *read()*, and
- 4689 the ability to accept an *aio_write()* and an *aio_cancel()* for a device that accepts a *write()*.
- 4690 This does not imply that the operation is asynchronous.

4691 **Standardization Issues**

4692 The following issues are addressed by the standardization of asynchronous I/O:

- 4693
- Rationale for New Interface

4694 Non-blocking I/O does not satisfy the needs of either realtime or high-performance
 4695 computing models; these models require that a process overlap program execution and
 4696 I/O processing. Realtime applications will often make use of direct I/O to or from the
 4697 address space of the process, or require synchronized (unbuffered) I/O; they also require
 4698 the ability to overlap this I/O with other computation. In addition, asynchronous I/O
 4699 allows an application to keep a device busy at all times, possibly achieving greater
 4700 throughput. Supercomputing and database architectures will often have specialized
 4701 hardware that can provide true asynchrony underlying the logical asynchrony provided
 4702 by this interface. In addition, asynchronous I/O should be supported by all types of files
 4703 and devices in the same manner.

- 4704
- Effect of Buffering

4705 If asynchronous I/O is performed on a file that is buffered prior to being actually written
 4706 to the device, it is possible that asynchronous I/O will offer no performance advantage
 4707 over normal I/O; the cycles *stolen* to perform the asynchronous I/O will be taken away
 4708 from the running process and the I/O will occur at interrupt time. This potential lack of
 4709 gain in performance in no way obviates the need for asynchronous I/O by realtime
 4710 applications, which very often will use specialized hardware support, multiple processors,
 4711 and/or unbuffered, synchronized I/O.

4712 **B.2.8.3 Memory Management**

4713 All memory management and shared memory definitions are located in the `<sys/mman.h>`
 4714 header. This is for alignment with historical practice.

4715 IEEE Std 1003.1-2001/Cor 1-2002, item XSH/TC1/D6/7 is applied, correcting the shading and
 4716 margin markers in the introduction to Section 2.8.3.1.

4717
4718
4719
4720
4721
4722
4723
4724
4725
4726
4727
4728
4729
4730
4731
4732
4733
4734
4735
4736
4737
4738
4739
4740
4741
4742
4743
4744
4745
4746
4747
4748
4749
4750
4751
4752
4753
4754
4755
4756
4757
4758
4759
4760
4761
4762
4763
4764
4765

Memory Locking Functions

This portion of the rationale presents models, requirements, and standardization issues relevant to process memory locking.

- Models

Realtime systems that conform to IEEE Std 1003.1-200x are expected (and desired) to be supported on systems with demand-paged virtual memory management, non-paged swapping memory management, and physical memory systems with no memory management hardware. The general case, however, is the demand-paged, virtual memory system with each POSIX process running in a virtual address space. Note that this includes architectures where each process resides in its own virtual address space and architectures where the address space of each process is only a portion of a larger global virtual address space.

The concept of memory locking is introduced to eliminate the indeterminacy introduced by paging and swapping, and to support an upper bound on the time required to access the memory mapped into the address space of a process. Ideally, this upper bound will be the same as the time required for the processor to access “main memory”, including any address translation and cache miss overheads. But some implementations—primarily on mainframes—will not actually force locked pages to be loaded and held resident in main memory. Rather, they will handle locked pages so that accesses to these pages will meet the performance metrics for locked process memory in the implementation. Also, although it is not, for example, the intention that this interface, as specified, be used to lock process memory into “cache”, it is conceivable that an implementation could support a large static RAM memory and define this as “main memory” and use a large[r] dynamic RAM as “backing store”. These interfaces could then be interpreted as supporting the locking of process memory into the static RAM. Support for multiple levels of backing store would require extensions to these interfaces.

Implementations may also use memory locking to guarantee a fixed translation between virtual and physical addresses where such is beneficial to improving determinacy for direct-to/from-process input/output. IEEE Std 1003.1-200x does not guarantee to the application that the virtual-to-physical address translations, if such exist, are fixed, because such behavior would not be implementable on all architectures on which implementations of IEEE Std 1003.1-200x are expected. But IEEE Std 1003.1-200x does mandate that an implementation define, for the benefit of potential users, whether or not locking guarantees fixed translations.

Memory locking is defined with respect to the address space of a process. Only the pages mapped into the address space of a process may be locked by the process, and when the pages are no longer mapped into the address space—for whatever reason—the locks established with respect to that address space are removed. Shared memory areas warrant special mention, as they may be mapped into more than one address space or mapped more than once into the address space of a process; locks may be established on pages within these areas with respect to several of these mappings. In such a case, the lock state of the underlying physical pages is the logical OR of the lock state with respect to each of the mappings. Only when all such locks have been removed are the shared pages considered unlocked.

In recognition of the page granularity of Memory Management Units (MMU), and in order to support locking of ranges of address space, memory locking is defined in terms of “page” granularity. That is, for the interfaces that support an address and size specification for the region to be locked, the address must be on a page boundary, and all pages mapped by the specified range are locked, if valid. This means that the length is implicitly rounded

4766 up to a multiple of the page size. The page size is implementation-defined and is available
4767 to applications as a compile-time symbolic constant or at runtime via *sysconf()*.

4768 A “real memory” POSIX.1b implementation that has no MMU could elect not to support
4769 these interfaces, returning [ENOSYS]. But an application could easily interpret this as
4770 meaning that the implementation would unconditionally page or swap the application
4771 when such is not the case. It is the intention of IEEE Std 1003.1-200x that such a system
4772 could define these interfaces as “NO-OPs”, returning success without actually performing
4773 any function except for mandated argument checking.

4774 • Requirements

4775 For realtime applications, memory locking is generally considered to be required as part of
4776 application initialization. This locking is performed after an application has been loaded
4777 (that is, *exec'd*) and the program remains locked for its entire lifetime. But to support
4778 applications that undergo major mode changes where, in one mode, locking is required,
4779 but in another it is not, the specified interfaces allow repeated locking and unlocking of
4780 memory within the lifetime of a process.

4781 When a realtime application locks its address space, it should not be necessary for the
4782 application to then “touch” all of the pages in the address space to guarantee that they are
4783 resident or else suffer potential paging delays the first time the page is referenced. Thus,
4784 IEEE Std 1003.1-200x requires that the pages locked by the specified interfaces be resident
4785 when the locking functions return successfully.

4786 Many architectures support system-managed stacks that grow automatically when the
4787 current extent of the stack is exceeded. A realtime application has a requirement to be able
4788 to “preallocate” sufficient stack space and lock it down so that it will not suffer page faults
4789 to grow the stack during critical realtime operation. There was no consensus on a portable
4790 way to specify how much stack space is needed, so IEEE Std 1003.1-200x supports no
4791 specific interface for preallocating stack space. But an application can portably lock down a
4792 specific amount of stack space by specifying *MCL_FUTURE* in a call to *mlockall()* and then
4793 calling a dummy function that declares an automatic array of the desired size.

4794 Memory locking for realtime applications is also generally considered to be an “all or
4795 nothing” proposition. That is, the entire process, or none, is locked down. But, for
4796 applications that have well-defined sections that need to be locked and others that do not,
4797 IEEE Std 1003.1-200x supports an optional set of interfaces to lock or unlock a range of
4798 process addresses. Reasons for locking down a specific range include:

- 4799 — An asynchronous event handler function that must respond to external events in a
4800 deterministic manner such that page faults cannot be tolerated
- 4801 — An input/output “buffer” area that is the target for direct-to-process I/O, and the
4802 overhead of implicit locking and unlocking for each I/O call cannot be tolerated

4803 Finally, locking is generally viewed as an “application-wide” function. That is, the
4804 application is globally aware of which regions are locked and which are not over time. This
4805 is in contrast to a function that is used temporarily within a “third party” library routine
4806 whose function is unknown to the application, and therefore must have no “side effects”.
4807 The specified interfaces, therefore, do not support “lock stacking” or “lock nesting” within
4808 a process. But, for pages that are shared between processes or mapped more than once
4809 into a process address space, “lock stacking” is essentially mandated by the requirement
4810 that unlocking of pages that are mapped by more than one process or more than once by
4811 the same process does not affect locks established on the other mappings.

4812 There was some support for “lock stacking” so that locking could be transparently used in
4813 functions or opaque modules. But the consensus was not to burden all implementations

4814 with lock stacking (and reference counting), and an implementation option was proposed.
 4815 There were strong objections to the option because applications would have to support
 4816 both options in order to remain portable. The consensus was to eliminate lock stacking
 4817 altogether, primarily through overwhelming support for the System V “m[un]lock[all]”
 4818 interface on which IEEE Std 1003.1-200x is now based.

4819 Locks are not inherited across *fork()*s because some implementations implement *fork()*
 4820 by creating new address spaces for the child. In such an implementation, requiring locks to be
 4821 inherited would lead to new situations in which a fork would fail due to the inability of
 4822 the system to lock sufficient memory to lock both the parent and the child. The consensus
 4823 was that there was no benefit to such inheritance. Note that this does not mean that locks
 4824 are removed when, for instance, a thread is created in the same address space.

4825 Similarly, locks are not inherited across *exec* because some implementations implement *exec*
 4826 by unmapping all of the pages in the address space (which, by definition, removes the
 4827 locks on these pages), and maps in pages of the *exec*'d image. In such an implementation,
 4828 requiring locks to be inherited would lead to new situations in which *exec* would fail.
 4829 Reporting this failure would be very cumbersome to detect in time to report to the calling
 4830 process, and no appropriate mechanism exists for informing the *exec*'d process of its status.

4831 It was determined that, if the newly loaded application required locking, it was the
 4832 responsibility of that application to establish the locks. This is also in keeping with the
 4833 general view that it is the responsibility of the application to be aware of all locks that are
 4834 established.

4835 There was one request to allow (not mandate) locks to be inherited across *fork()*, and a
 4836 request for a flag, *MCL_INHERIT*, that would specify inheritance of memory locks across
 4837 *exec*s. Given the difficulties raised by this and the general lack of support for the feature in
 4838 IEEE Std 1003.1-200x, it was not added. IEEE Std 1003.1-200x does not preclude an
 4839 implementation from providing this feature for administrative purposes, such as a “run”
 4840 command that will lock down and execute a specified application. Additionally, the
 4841 rationale for the objection equated *fork()* with creating a thread in the address space.
 4842 IEEE Std 1003.1-200x does not mandate releasing locks when creating additional threads in
 4843 an existing process.

4844 • Standardization Issues

4845 One goal of IEEE Std 1003.1-200x is to define a set of primitives that provide the necessary
 4846 functionality for realtime applications, with consideration for the needs of other
 4847 application domains where such were identified, which is based to the extent possible on
 4848 existing industry practice.

4849 The Memory Locking option is required by many realtime applications to tune
 4850 performance. Such a facility is accomplished by placing constraints on the virtual memory
 4851 system to limit paging of time of the process or of critical sections of the process. This
 4852 facility should not be used by most non-realtime applications.

4853 Optional features provided in IEEE Std 1003.1-200x allow applications to lock selected
 4854 address ranges with the caveat that the process is responsible for being aware of the page
 4855 granularity of locking and the unnested nature of the locks.

4856
4857
4858
4859
4860
4861
4862
4863
4864
4865
4866
4867
4868
4869
4870
4871
4872
4873
4874
4875
4876
4877
4878
4879
4880
4881
4882
4883
4884
4885
4886
4887
4888
4889
4890
4891
4892
4893
4894
4895
4896
4897
4898
4899
4900
4901
4902

Mapped Files Functions

The memory mapped files functionality provides a mechanism that allows a process to access files by directly incorporating file data into its address space. Once a file is “mapped” into a process address space, the data can be manipulated by instructions as memory. The use of mapped files can significantly reduce I/O data movement since file data does not have to be copied into process data buffers as in *read()* and *write()*. If more than one process maps a file, its contents are shared among them. This provides a low overhead mechanism by which processes can synchronize and communicate.

- Historical Perspective

Realtime applications have historically been implemented using a collection of cooperating processes or tasks. In early systems, these processes ran on bare hardware (that is, without an operating system) with no memory relocation or protection. The application paradigms that arose from this environment involve the sharing of data between the processes.

When realtime systems were implemented on top of vendor-supplied operating systems, the paradigm or performance benefits of direct access to data by multiple processes was still deemed necessary. As a result, operating systems that claim to support realtime applications must support the shared memory paradigm.

Additionally, a number of realtime systems provide the ability to map specific sections of the physical address space into the address space of a process. This ability is required if an application is to obtain direct access to memory locations that have specific properties (for example, refresh buffers or display devices, dual ported memory locations, DMA target locations). The use of this ability is common enough to warrant some degree of standardization of its interface. This ability overlaps the general paradigm of shared memory in that, in both instances, common global objects are made addressable by individual processes or tasks.

Finally, a number of systems also provide the ability to map process addresses to files. This provides both a general means of sharing persistent objects, and using files in a manner that optimizes memory and swapping space usage.

Simple shared memory is clearly a special case of the more general file mapping capability. In addition, there is relatively widespread agreement and implementation of the file mapping interface. In these systems, many different types of objects can be mapped (for example, files, memory, devices, and so on) using the same mapping interfaces. This approach both minimizes interface proliferation and maximizes the generality of programs using the mapping interfaces.

- Memory Mapped Files Usage

A memory object can be concurrently mapped into the address space of one or more processes. The *mmap()* and *munmap()* functions allow a process to manipulate their address space by mapping portions of memory objects into it and removing them from it. When multiple processes map the same memory object, they can share access to the underlying data. Implementations may restrict the size and alignment of mappings to be on *page*-size boundaries. The page size, in bytes, is the value of the system-configurable variable {PAGESIZE}, typically accessed by calling *sysconf()* with a *name* argument of *_SC_PAGESIZE*. If an implementation has no restrictions on size or alignment, it may specify a 1-byte page size.

To map memory, a process first opens a memory object. The *ftruncate()* function can be used to contract or extend the size of the memory object even when the object is currently mapped. If the memory object is extended, the contents of the extended areas are zeros.

4903 After opening a memory object, the application maps the object into its address space
 4904 using the *mmap()* function call. Once a mapping has been established, it remains mapped
 4905 until unmapped with *munmap()*, even if the memory object is closed. The *mprotect()*
 4906 function can be used to change the memory protections initially established by *mmap()*.

4907 A *close()* of the file descriptor, while invalidating the file descriptor itself, does not unmap
 4908 any mappings established for the memory object. The address space, including all mapped
 4909 regions, is inherited on *fork()*. The entire address space is unmapped on process
 4910 termination or by successful calls to any of the *exec* family of functions.

4911 The *msync()* function is used to force mapped file data to permanent storage.

4912 • Effects on Other Functions

4913 With memory mapped files, the operation of the *open()*, *creat()*, and *unlink()* functions are
 4914 a natural result of using the file system name space to map the global names for memory
 4915 objects.

4916 The *ftruncate()* function can be used to set the length of a sharable memory object.

4917 The meaning of *stat()* fields other than the size and protection information is undefined on
 4918 implementations where memory objects are not implemented using regular files. When
 4919 regular files are used, the times reflect when the implementation updated the file image of
 4920 the data, not when a process updated the data in memory.

4921 The operations of *fdopen()*, *write()*, *read()*, and *lseek()* were made unspecified for objects
 4922 opened with *shm_open()*, so that implementations that did not implement memory objects
 4923 as regular files would not have to support the operation of these functions on shared
 4924 memory objects.

4925 The behavior of memory objects with respect to *close()*, *dup()*, *dup2()*, *open()*, *close()*,
 4926 *fork()*, *_exit()*, and the *exec* family of functions is the same as the behavior of the existing
 4927 practice of the *mmap()* function.

4928 A memory object can still be referenced after a close. That is, any mappings made to the
 4929 file are still in effect, and reads and writes that are made to those mappings are still valid
 4930 and are shared with other processes that have the same mapping. Likewise, the memory
 4931 object can still be used if any references remain after its name(s) have been deleted. Any
 4932 references that remain after a close must not appear to the application as file descriptors.

4933 This is existing practice for *mmap()* and *close()*. In addition, there are already mappings
 4934 present (text, data, stack) that do not have open file descriptors. The text mapping in
 4935 particular is considered a reference to the file containing the text. The desire was to treat all
 4936 mappings by the process uniformly. Also, many modern implementations use *mmap()* to
 4937 implement shared libraries, and it would not be desirable to keep file descriptors for each
 4938 of the many libraries an application can use. It was felt there were many other existing
 4939 programs that used this behavior to free a file descriptor, and thus IEEE Std 1003.1-200x
 4940 could not forbid it and still claim to be using existing practice.

4941 For implementations that implement memory objects using memory only, memory objects
 4942 will retain the memory allocated to the file after the last close and will use that same
 4943 memory on the next open. Note that closing the memory object is not the same as deleting
 4944 the name, since the memory object is still defined in the memory object name space.

4945 The locks of *fcntl()* do not block any read or write operation, including read or write access
 4946 to shared memory or mapped files. In addition, implementations that only support shared
 4947 memory objects should not be required to implement record locks. The reference to *fcntl()*
 4948 is added to make this point explicitly. The other *fcntl()* commands are useful with shared
 4949 memory objects.

4950 The size of pages that mapping hardware may be able to support may be a configurable
 4951 value, or it may change based on hardware implementations. The addition of the
 4952 `_SC_PAGESIZE` parameter to the `sysconf()` function is provided for determining the
 4953 mapping page size at runtime.

4954 Shared Memory Functions

4955 Implementations may support the Shared Memory Objects option independently of memory
 4956 mapped files. Shared memory objects are named regions of storage that may be independent of
 4957 the file system and can be mapped into the address space of one or more processes to allow
 4958 them to share the associated memory.

4959 • Requirements

4960 Shared memory is used to share data among several processes, each potentially running at
 4961 different priority levels, responding to different inputs, or performing separate tasks.
 4962 Shared memory is not just simply providing common access to data, it is providing the
 4963 fastest possible communication between the processes. With one memory write operation,
 4964 a process can pass information to as many processes as have the memory region mapped.

4965 As a result, shared memory provides a mechanism that can be used for all other
 4966 interprocess communication facilities. It may also be used by an application for
 4967 implementing more sophisticated mechanisms than semaphores and message queues.

4968 The need for a shared memory interface is obvious for virtual memory systems, where the
 4969 operating system is directly preventing processes from accessing each other's data.
 4970 However, in unprotected systems, such as those found in some embedded controllers, a
 4971 shared memory interface is needed to provide a portable mechanism to allocate a region of
 4972 memory to be shared and then to communicate the address of that region to other
 4973 processes.

4974 This, then, provides the minimum functionality that a shared memory interface must have
 4975 in order to support realtime applications: to allocate and name an object to be mapped into
 4976 memory for potential sharing (`open()` or `shm_open()`), and to make the memory object
 4977 available within the address space of a process (`mmap()`). To complete the interface, a
 4978 mechanism to release the claim of a process on a shared memory object (`munmap()`) is also
 4979 needed, as well as a mechanism for deleting the name of a sharable object that was
 4980 previously created (`unlink()` or `shm_unlink()`).

4981 After a mapping has been established, an implementation should not have to provide
 4982 services to maintain that mapping. All memory writes into that area will appear
 4983 immediately in the memory mapping of that region by any other processes.

4984 Thus, requirements include:

- 4985 — Support creation of sharable memory objects and the mapping of these objects into
 4986 the address space of a process.
- 4987 — Sharable memory objects should be accessed by global names accessible from all
 4988 processes.
- 4989 — Support the mapping of specific sections of physical address space (such as a
 4990 memory mapped device) into the address space of a process. This should not be
 4991 done by the process specifying the actual address, but again by an implementation-
 4992 defined global name (such as a special device name) dedicated to this purpose.
- 4993 — Support the mapping of discrete portions of these memory objects.

- 4994 — Support for minimum hardware configurations that contain no physical media on
- 4995 which to store shared memory contents permanently.
- 4996 — The ability to preallocate the entire shared memory region so that minimum
- 4997 hardware configurations without virtual memory support can guarantee contiguous
- 4998 space.
- 4999 — The maximizing of performance by not requiring functionality that would require
- 5000 implementation interaction above creating the shared memory area and returning
- 5001 the mapping.

5002 Note that the above requirements do not preclude:

- 5003 — The sharable memory object from being implemented using actual files on an actual
- 5004 file system.
- 5005 — The global name that is accessible from all processes being restricted to a file system
- 5006 area that is dedicated to handling shared memory.
- 5007 — An implementation not providing implementation-defined global names for the
- 5008 purpose of physical address mapping.

5009 • Shared Memory Objects Usage

5010 If the Shared Memory Objects option is supported, a shared memory object may be
 5011 created, or opened if it already exists, with the *shm_open()* function. If the shared memory
 5012 object is created, it has a length of zero. The *ftruncate()* function can be used to set the size
 5013 of the shared memory object after creation. The *shm_unlink()* function removes the name
 5014 for a shared memory object created by *shm_open()*.

5015 • Shared Memory Overview

5016 The shared memory facility defined by IEEE Std 1003.1-200x usually results in memory
 5017 locations being added to the address space of the process. The implementation returns the
 5018 address of the new space to the application by means of a pointer. This works well in
 5019 languages like C. However, in languages without pointer types it will not work. In the
 5020 bindings for such a language, either a special COMMON section will need to be defined
 5021 (which is unlikely), or the binding will have to allow existing structures to be mapped. The
 5022 implementation will likely have to place restrictions on the size and alignment of such
 5023 structures or will have to map a suitable region of the address space of the process into the
 5024 memory object, and thus into other processes. These are issues for that particular language
 5025 binding. For IEEE Std 1003.1-200x, however, the practice will not be forbidden, merely
 5026 undefined.

5027 Two potentially different name spaces are used for naming objects that may be mapped
 5028 into process address spaces. When using memory mapped files, files may be accessed via
 5029 *open()*. When the Shared Memory Objects option is supported, sharable memory objects
 5030 that might not be files may be accessed via the *shm_open()* function. These operations are
 5031 not mutually-exclusive.

5032 Some implementations supporting the Shared Memory Objects option may choose to
 5033 implement the shared memory object name space as part of the file system name space.
 5034 There are several reasons for this:

- 5035 — It allows applications to prevent name conflicts by use of the directory structure.
- 5036 — It uses an existing mechanism for accessing global objects and prevents the creation
- 5037 of a new mechanism for naming global objects.

5038 In such implementations, memory objects can be implemented using regular files, if that is

5039 what the implementation chooses. The *shm_open()* function can be implemented as an
 5040 *open()* call in a fixed directory followed by a call to *fcntl()* to set *FD_CLOEXEC*. The
 5041 *shm_unlink()* function can be implemented as an *unlink()* call.

5042 On the other hand, it is also expected that small embedded systems that support the
 5043 Shared Memory Objects option may wish to implement shared memory without having
 5044 any file systems present. In this case, the implementations may choose to use a simple
 5045 string valued name space for shared memory regions. The *shm_open()* function permits
 5046 either type of implementation.

5047 Some implementations have hardware that supports protection of mapped data from
 5048 certain classes of access and some do not. Systems that supply this functionality support
 5049 the memory protection functionality.

5050 Some implementations restrict size, alignment, and protections to be on *page-size*
 5051 boundaries. If an implementation has no restrictions on size or alignment, it may specify a
 5052 1-byte page size. Applications on implementations that do support larger pages must be
 5053 cognizant of the page size since this is the alignment and protection boundary.

5054 Simple embedded implementations may have a 1-byte page size and only support the
 5055 Shared Memory Objects option. This provides simple shared memory between processes
 5056 without requiring mapping hardware.

5057 IEEE Std 1003.1-200x specifically allows a memory object to remain referenced after a close
 5058 because that is existing practice for the *mmap()* function.

5059 Typed Memory Functions

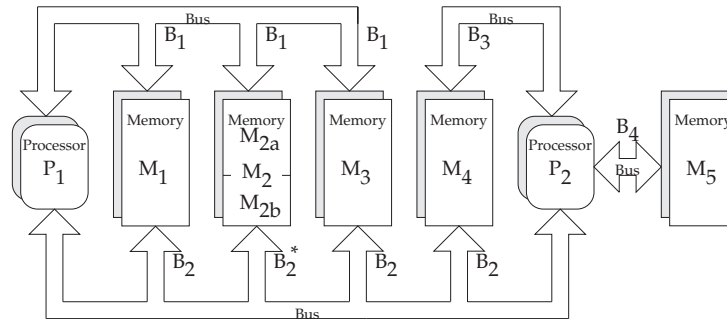
5060 Implementations may support the Typed Memory Objects option without supporting either the
 5061 Shared Memory option or memory mapped files. Typed memory objects are pools of specialized
 5062 storage, different from the main memory resource normally used by a processor to hold code
 5063 and data, that can be mapped into the address space of one or more processes.

- 5064 • Model

5065 Realtime systems conforming to one of the POSIX.13 realtime profiles are expected (and
 5066 desired) to be supported on systems with more than one type or pool of memory (for
 5067 example, SRAM, DRAM, ROM, EPROM, EEPROM), where each type or pool of memory
 5068 may be accessible by one or more processors via one or more buses (ports). Memory
 5069 mapped files, shared memory objects, and the language-specific storage allocation
 5070 operators (*malloc()* for the ISO C standard, *new* for ISO Ada) fail to provide application
 5071 program interfaces versatile enough to allow applications to control their utilization of
 5072 such diverse memory resources. The typed memory interfaces *posix_typed_mem_open()*,
 5073 *posix_mem_offset()*, *posix_typed_mem_get_info()*, *mmap()*, and *munmap()* defined herein
 5074 support the model of typed memory described below.

5075 For purposes of this model, a system comprises several processors (for example, P_1 and
 5076 P_2), several physical memory pools (for example, M_1 , M_2 , M_{2a} , M_{2b} , M_3 , M_4 , and M_5), and
 5077 several buses or “ports” (for example, B_1 , B_2 , B_3 , and B_4) interconnecting the various
 5078 processors and memory pools in some system-specific way. Notice that some memory
 5079 pools may be contained in others (for example, M_{2a} and M_{2b} are contained in M_2).

5080 **Figure B-1** shows an example of such a model. In a system like this, an application should
 5081 be able to perform the following operations:



* All addresses in pool M₂ (comprising pools M_{2a} and M_{2b}) accessible via port B₁.
 Addresses in pool M_{2b} are also accessible via port B₂.
 Addresses in pool M_{2a} are *not* accessible via port B₂.

5082

Figure B-1 Example of a System with Typed Memory

5083

— Typed Memory Allocation

5084

5085

5086

5087

5088

5089

An application should be able to allocate memory dynamically from the desired pool using the desired bus, and map it into the address space of a process. For example, processor P₁ can allocate some portion of memory pool M₁ through port B₁, treating all unmapped subareas of M₁ as a heap-storage resource from which memory may be allocated. This portion of memory is mapped into address space of the process, and subsequently deallocated when unmapped from all processes.

5090

— Using the Same Storage Region from Different Busses

5091

5092

5093

5094

5095

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 *mmap()*), to allow it to be accessed from another bus. For example, processor P₁ may wish to access the same region of memory pool M_{2b} both through ports B₁ and B₂.

5096

— Sharing Typed Memory Regions

5097

5098

5099

5100

5101

5102

5103

5104

5105

5106

5107

5108

5109

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 P₁ may want to share a region of memory pool M₄ with processor P₂, and they may be required to use buses B₂ and B₃, 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.

5110

— Contiguous Allocation

5111

5112

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

5113 solved by requesting contiguous allocation. For example, a process in P_1 can allocate
 5114 10 Kbytes of physically contiguous memory from M_3-B_1 , and obtain the offset (within
 5115 pool M_3) of this block of memory. Then, it can pass this offset (and the length) to a
 5116 process in P_2 using some interprocess communication mechanism. The second
 5117 process can map the same block of memory by using the offset transferred and
 5118 specifying M_3-B_2 .

5119 — Unallocated Mapping

5120 Any subarea of a memory pool that is mapped to a process, either as the result of an
 5121 allocation request or an explicit mapping, is normally unavailable for allocation.
 5122 Special processes such as debuggers, however, may need to map large areas of a
 5123 typed memory pool, yet leave those areas available for allocation.

5124 Typed memory allocation and mapping has to coexist with storage allocation operators
 5125 like *malloc()*, but systems are free to choose how to implement this coexistence. For
 5126 example, it may be system configuration-dependent if all available system memory is
 5127 made part of one of the typed memory pools or if some part will be restricted to
 5128 conventional allocation operators. Equally system configuration-dependent may be the
 5129 availability of operators like *malloc()* to allocate storage from certain typed memory pools.
 5130 It is not excluded to configure a system such that a given named pool, P_1 , is in turn split
 5131 into non-overlapping named subpools. For example, M_1-B_1 , M_2-B_1 , and M_3-B_1 could also be
 5132 accessed as one common pool $M_{123}-B_1$. A call to *malloc()* on P_1 could work on such a larger
 5133 pool while full optimization of memory usage by P_1 would require typed memory
 5134 allocation at the subpool level.

5135 • Existing Practice

5136 OS-9 provides for the naming (numbering) and prioritization of memory types by a system
 5137 administrator. It then provides APIs to request memory allocation of typed (colored)
 5138 memory by number, and to generate a bus address from a mapped memory address
 5139 (translate). When requesting colored memory, the user can specify type 0 to signify
 5140 allocation from the first available type in priority order.

5141 HP-RT presents interfaces to map different kinds of storage regions that are visible through
 5142 a VME bus, although it does not provide allocation operations. It also provides functions
 5143 to perform address translation between VME addresses and virtual addresses. It represents
 5144 a VME-bus unique solution to the general problem.

5145 The PSOS approach is similar (that is, based on a pre-established mapping of bus address
 5146 ranges to specific memories) with a concept of segments and regions (regions dynamically
 5147 allocated from a heap which is a special segment). Therefore, PSOS does not fully address
 5148 the general allocation problem either. PSOS does not have a “process”-based model, but
 5149 more of a “thread”-only-based model of multi-tasking. So mapping to a process address
 5150 space is not an issue.

5151 QNX uses the System V approach of opening specially named devices (shared memory
 5152 segments) and using *mmap()* to then gain access from the process. They do not address
 5153 allocation directly, but once typed shared memory can be mapped, an “allocation
 5154 manager” process could be written to handle requests for allocation.

5155 The System V approach also included allocation, implemented by opening yet other
 5156 special “devices” which allocate, rather than appearing as a whole memory object.

5157 The Orkid realtime kernel interface definition has operations to manage memory “regions”
 5158 and “pools”, which are areas of memory that may reflect the differing physical nature of
 5159 the memory. Operations to allocate memory from these regions and pools are also
 5160 provided.

- 5161
- Requirements
- 5162 Existing practice in SVID-derived UNIX systems relies on functionality similar to *mmap()*
5163 and its related interfaces to achieve mapping and allocation of typed memory. However,
5164 the issue of sharing typed memory (allocated or mapped) and the complication of multiple
5165 ports are not addressed in any consistent way by existing UNIX system practice. Part of
5166 this functionality is existing practice in specialized realtime operating systems. In order to
5167 solidify the capabilities implied by the model above, the following requirements are
5168 imposed on the interface:
- Identification of Typed Memory Pools and Ports

5169

5170 All processes (running in all processors) in the system are able to identify a particular
5171 (system configured) typed memory pool accessed through a particular (system
5172 configured) port by a name. That name is a member of a name space common to all
5173 these processes, but need not be the same name space as that containing ordinary
5174 filenames. The association between memory pools/ports and corresponding names
5175 is typically established when the system is configured. The “open” operation for
5176 typed memory objects should be distinct from the *open()* function, for consistency
5177 with other similar services, but implementable on top of *open()*. This implies that the
5178 handle for a typed memory object will be a file descriptor.
 - Allocation and Mapping of Typed Memory

5179

5180 Once a typed memory object has been identified by a process, it is possible to both
5181 map user-selected subareas of that object into process address space and to map
5182 system-selected (that is, dynamically allocated) subareas of that object, with user-
5183 specified length, into process address space. It is also possible to determine the
5184 maximum length of memory allocation that may be requested from a given typed
5185 memory object.
 - Sharing Typed Memory

5186

5187 Two or more processes are able to share portions of typed memory, either user-
5188 selected or dynamically allocated. This requirement applies also to dynamically
5189 allocated regions of memory that are composed of several non-contiguous pieces.
 - Contiguous Allocation

5190

5191 For dynamic allocation, it is the user’s option whether the system is required to
5192 allocate a contiguous subarea within the typed memory object, or whether it is
5193 permitted to allocate discontinuous fragments which appear contiguous in the
5194 process mapping. Contiguous allocation simplifies the process of sharing allocated
5195 typed memory, while discontinuous allocation allows for potentially better recovery
5196 of deallocated typed memory.
 - Accessing Typed Memory Through Different Ports

5197

5198 Once a subarea of a typed memory object has been mapped, it is possible to
5199 determine the location and length corresponding to a user-selected portion of that
5200 object within the memory pool. This location and length can then be used to remap
5201 that portion of memory for access from another port. If the referenced portion of
5202 typed memory was allocated discontinuously, the length thus determined may be
5203 shorter than anticipated, and the user code must adapt to the value returned.
 - Deallocation

5204

5205 When a previously mapped subarea of typed memory is no longer mapped by any
5206 process in the system—as a result of a call or calls to *munmap()*—that subarea

5207 becomes potentially reusable for dynamic allocation; actual reuse of the subarea is a
5208 function of the dynamic typed memory allocation policy.

5209 — Unallocated Mapping

5210 It must be possible to map user-selected subareas of a typed memory object without
5211 marking that subarea as unavailable for allocation. This option is not the default
5212 behavior, and requires appropriate privilege.

5213 • Scenario

5214 The following scenario will serve to clarify the use of the typed memory interfaces.

5215 Process A running on P_1 (see Figure B-1 (on page 122)) wants to allocate some memory
5216 from memory pool M_2 , and it wants to share this portion of memory with process B
5217 running on P_2 . Since P_2 only has access to the lower part of M_2 , both processes will use the
5218 memory pool named M_{2b} which is the part of M_2 that is accessible both from P_1 and P_2 . The
5219 operations that both processes need to perform are shown below:

5220 — Allocating Typed Memory

5221 Process A calls `posix_typed_mem_open()` with the name `/typed.m2b-b1` and a `tflag` of
5222 `POSIX_TYPED_MEM_ALLOCATE` to get a file descriptor usable for allocating from
5223 pool M_{2b} accessed through port B_1 . It then calls `mmap()` with this file descriptor
5224 requesting a length of 4096 bytes. The system allocates two discontinuous blocks of
5225 sizes 1024 and 3072 bytes within M_{2b} . The `mmap()` function returns a pointer to a
5226 4096-byte array in process A's logical address space, mapping the allocated blocks
5227 contiguously. Process A can then utilize the array, and store data in it.

5228 — Determining the Location of the Allocated Blocks

5229 Process A can determine the lengths and offsets (relative to M_{2b}) of the two blocks
5230 allocated, by using the following procedure: First, process A calls `posix_mem_offset()`
5231 with the address of the first element of the array and length 4096. Upon return, the
5232 offset and length (1024 bytes) of the first block are returned. A second call to
5233 `posix_mem_offset()` is then made using the address of the first element of the array
5234 plus 1024 (the length of the first block), and a new length of 4096-1024. If there were
5235 more fragments allocated, this procedure could have been continued within a loop
5236 until the offsets and lengths of all the blocks were obtained. Notice that this relatively
5237 complex procedure can be avoided if contiguous allocation is requested (by opening
5238 the typed memory object with the `tflag`
5239 `POSIX_TYPED_MEM_ALLOCATE_CONTIG`).

5240 — Sharing Data Across Processes

5241 Process A passes the two offset values and lengths obtained from the
5242 `posix_mem_offset()` calls to process B running on P_2 , via some form of interprocess
5243 communication. Process B can gain access to process A's data by calling
5244 `posix_typed_mem_open()` with the name `/typed.m2b-b2` and a `tflag` of zero, then using
5245 two `mmap()` calls on the resulting file descriptor to map the two subareas of that
5246 typed memory object to its own address space.

5247 • Rationale for no `mem_alloc()` and `mem_free()`

5248 The standard developers had originally proposed a pair of new flags to `mmap()` which,
5249 when applied to a typed memory object descriptor, would cause `mmap()` to allocate
5250 dynamically from an unallocated and unmapped area of the typed memory object.
5251 Deallocation was similarly accomplished through the use of `munmap()`. This was rejected
5252 by the ballot group because it excessively complicated the (already rather complex)

5253 *mmap()* interface and introduced semantics useful only for typed memory, to a function
 5254 which must also map shared memory and files. They felt that a memory allocator should
 5255 be built on top of *mmap()* instead of being incorporated within the same interface, much as
 5256 the ISO C standard libraries build *malloc()* on top of the virtual memory mapping
 5257 functions *brk()* and *sbrk()*. This would eliminate the complicated semantics involved with
 5258 unmapping only part of an allocated block of typed memory.

5259 To attempt to achieve ballot group consensus, typed memory allocation and deallocation
 5260 was first migrated from *mmap()* and *munmap()* to a pair of complementary functions
 5261 modeled on the ISO C standard *malloc()* and *free()*. The *mem_alloc()* function specified
 5262 explicitly the typed memory object (typed memory pool/access port) from which
 5263 allocation takes place, unlike *malloc()* where the memory pool and port are unspecified.
 5264 The *mem_free()* function handled deallocation. These new semantics still met all of the
 5265 requirements detailed above without modifying the behavior of *mmap()* except to allow it
 5266 to map specified areas of typed memory objects. An implementation would have been free
 5267 to implement *mem_alloc()* and *mem_free()* over *mmap()*, through *mmap()*, or independently
 5268 but cooperating with *mmap()*.

5269 The ballot group was queried to see if this was an acceptable alternative, and while there
 5270 was some agreement that it achieved the goal of removing the complicated semantics of
 5271 allocation from the *mmap()* interface, several balloters realized that it just created two
 5272 additional functions that behaved, in great part, like *mmap()*. These balloters proposed an
 5273 alternative which has been implemented here in place of a separate *mem_alloc()* and
 5274 *mem_free()*. This alternative is based on four specific suggestions:

- 5275 1. The *posix_typed_mem_open()* function should provide a flag which specifies
 5276 “allocate on *mmap()*” (otherwise, *mmap()* just maps the underlying object). This
 5277 allows things roughly similar to */dev/zero* versus */dev/swap*. Two such flags have
 5278 been implemented, one of which forces contiguous allocation.
- 5279 2. The *posix_mem_offset()* function is acceptable because it can be applied usefully to
 5280 mapped objects in general. It should return the file descriptor of the underlying
 5281 object.
- 5282 3. The *mem_get_info()* function in an earlier draft should be renamed
 5283 *posix_typed_mem_get_info()* because it is not generally applicable to memory objects.
 5284 It should probably return the file descriptor’s allocation attribute. The renaming of
 5285 the function has been implemented, but having it return a piece of information
 5286 which is readily known by an application without this function has been rejected.
 5287 Its whole purpose is to query the typed memory object for attributes that are not
 5288 user-specified, but determined by the implementation.
- 5289 4. There should be no separate *mem_alloc()* or *mem_free()* functions. Instead, using
 5290 *mmap()* on a typed memory object opened with an “allocate on *mmap()*” flag
 5291 should be used to force allocation. These are precisely the semantics defined in the
 5292 current draft.

5293 • Rationale for no Typed Memory Access Management

5294 The working group had originally defined an additional interface (and an additional kind
 5295 of object: typed memory master) to establish and dissolve mappings to typed memory on
 5296 behalf of devices or processors which were independent of the operating system and had
 5297 no inherent capability to directly establish mappings on their own. This was to have
 5298 provided functionality similar to device driver interfaces such as *physio()* and their
 5299 underlying bus-specific interfaces (for example, *mballoc()*) which serve to set up and break
 5300 down DMA pathways, and derive mapped addresses for use by hardware devices and
 5301 processor cards.

5302 The ballot group felt that this was beyond the scope of POSIX.1 and its amendments.
 5303 Furthermore, the removal of interrupt handling interfaces from a preceding amendment
 5304 (the IEEE Std 1003.1d-1999) during its balloting process renders these typed memory
 5305 access management interfaces an incomplete solution to portable device management from
 5306 a user process; it would be possible to initiate a device transfer to/from typed memory, but
 5307 impossible to handle the transfer-complete interrupt in a portable way.

5308 To achieve ballot group consensus, all references to typed memory access management
 5309 capabilities were removed. The concept of portable interfaces from a device driver to both
 5310 operating system and hardware is being addressed by the Uniform Driver Interface (UDI)
 5311 industry forum, with formal standardization deferred until proof of concept and industry-
 5312 wide acceptance and implementation.

5313 B.2.8.4 Process Scheduling

5314 IEEE PASC Interpretation 1003.1 #96 has been applied, adding the `pthread_setschedprio()`
 5315 function. This was added since previously there was no way for a thread to lower its own
 5316 priority without going to the tail of the threads list for its new priority. This capability is
 5317 necessary to bound the duration of priority inversion encountered by a thread.

5318 The following portion of the rationale presents models, requirements, and standardization
 5319 issues relevant to process scheduling; see also [Section B.2.9.4](#) (on page 166).

5320 In an operating system supporting multiple concurrent processes, the system determines the
 5321 order in which processes execute to meet implementation-defined goals. For time-sharing
 5322 systems, the goal is to enhance system throughput and promote fairness; the application is
 5323 provided with little or no control over this sequencing function. While this is acceptable and
 5324 desirable behavior in a time-sharing system, it is inappropriate in a realtime system; realtime
 5325 applications must specifically control the execution sequence of their concurrent processes in
 5326 order to meet externally defined response requirements.

5327 In IEEE Std 1003.1-200x, the control over process sequencing is provided using a concept of
 5328 scheduling policies. These policies, described in detail in this section, define the behavior of the
 5329 system whenever processor resources are to be allocated to competing processes. Only the
 5330 behavior of the policy is defined; conforming implementations are free to use any mechanism
 5331 desired to achieve the described behavior.

- 5332 • Models

5333 In an operating system supporting multiple concurrent processes, the system determines
 5334 the order in which processes execute and might force long-running processes to yield to
 5335 other processes at certain intervals. Typically, the scheduling code is executed whenever an
 5336 event occurs that might alter the process to be executed next.

5337 The simplest scheduling strategy is a “first-in, first-out” (FIFO) dispatcher. Whenever a
 5338 process becomes runnable, it is placed on the end of a ready list. The process at the front of
 5339 the ready list is executed until it exits or becomes blocked, at which point it is removed
 5340 from the list. This scheduling technique is also known as “run-to-completion” or “run-to-
 5341 block”.

5342 A natural extension to this scheduling technique is the assignment of a “non-migrating
 5343 priority” to each process. This policy differs from strict FIFO scheduling in only one
 5344 respect: whenever a process becomes runnable, it is placed at the end of the list of
 5345 processes runnable at that priority level. When selecting a process to run, the system
 5346 always selects the first process from the highest priority queue with a runnable process.
 5347 Thus, when a process becomes unblocked, it will preempt a running process of lower
 5348 priority without otherwise altering the ready list. Further, if a process elects to alter its
 5349 priority, it is removed from the ready list and reinserted, using its new priority, according

5350 to the policy above.

5351 While the above policy might be considered unfriendly in a time-sharing environment in
5352 which multiple users require more balanced resource allocation, it could be ideal in a
5353 realtime environment for several reasons. The most important of these is that it is
5354 deterministic: the highest-priority process is always run and, among processes of equal
5355 priority, the process that has been runnable for the longest time is executed first. Because of
5356 this determinism, cooperating processes can implement more complex scheduling simply
5357 by altering their priority. For instance, if processes at a single priority were to reschedule
5358 themselves at fixed time intervals, a time-slice policy would result.

5359 In a dedicated operating system in which all processes are well-behaved realtime
5360 applications, non-migrating priority scheduling is sufficient. However, many existing
5361 implementations provide for more complex scheduling policies.

5362 IEEE Std 1003.1-200x specifies a linear scheduling model. In this model, every process in
5363 the system has a priority. The system scheduler always dispatches a process that has the
5364 highest (generally the most time-critical) priority among all runnable processes in the
5365 system. As long as there is only one such process, the dispatching policy is trivial. When
5366 multiple processes of equal priority are eligible to run, they are ordered according to a
5367 strict run-to-completion (FIFO) policy.

5368 The priority is represented as a positive integer and is inherited from the parent process.
5369 For processes running under a fixed priority scheduling policy, the priority is never altered
5370 except by an explicit function call.

5371 It was determined arbitrarily that larger integers correspond to “higher priorities”.

5372 Certain implementations might impose restrictions on the priority ranges to which
5373 processes can be assigned. There also can be restrictions on the set of policies to which
5374 processes can be set.

5375 • Requirements

5376 Realtime processes require that scheduling be fast and deterministic, and that it guarantees
5377 to preempt lower priority processes.

5378 Thus, given the linear scheduling model, realtime processes require that they be run at a
5379 priority that is higher than other processes. Within this framework, realtime processes are
5380 free to yield execution resources to each other in a completely portable and
5381 implementation-defined manner.

5382 As there is a generally perceived requirement for processes at the same priority level to
5383 share processor resources more equitably, provisions are made by providing a scheduling
5384 policy (that is, SCHED_RR) intended to provide a timeslice-like facility.

5385 **Note:** The following topics assume that low numeric priority implies low scheduling criticality
5386 and *vice versa*.

5387 • Rationale for New Interface

5388 Realtime applications need to be able to determine when processes will run in relation to
5389 each other. It must be possible to guarantee that a critical process will run whenever it is
5390 runnable; that is, whenever it wants to for as long as it needs. SCHED_FIFO satisfies this
5391 requirement. Additionally, SCHED_RR was defined to meet a realtime requirement for a
5392 well-defined time-sharing policy for processes at the same priority.

5393 It would be possible to use the BSD *setpriority()* and *getpriority()* functions by redefining
5394 the meaning of the “nice” parameter according to the scheduling policy currently in use by
5395 the process. The System V *nice()* interface was felt to be undesirable for realtime because it

5396 specifies an adjustment to the “nice” value, rather than setting it to an explicit value.
 5397 Realtime applications will usually want to set priority to an explicit value. Also, System V
 5398 *nice()* does not allow for changing the priority of another process.

5399 With the POSIX.1b interfaces, the traditional “nice” value does not affect the SCHED_FIFO
 5400 or SCHED_RR scheduling policies. If a “nice” value is supported, it is implementation-
 5401 defined whether it affects the SCHED_OTHER policy.

5402 An important aspect of IEEE Std 1003.1-200x is the explicit description of the queuing and
 5403 preemption rules. It is critical, to achieve deterministic scheduling, that such rules be
 5404 stated clearly in IEEE Std 1003.1-200x.

5405 IEEE Std 1003.1-200x does not address the interaction between priority and swapping. The
 5406 issues involved with swapping and virtual memory paging are extremely implementation-
 5407 defined and would be nearly impossible to standardize at this point. The proposed
 5408 scheduling paradigm, however, fully describes the scheduling behavior of runnable
 5409 processes, of which one criterion is that the working set be resident in memory. Assuming
 5410 the existence of a portable interface for locking portions of a process in memory, paging
 5411 behavior need not affect the scheduling of realtime processes.

5412 IEEE Std 1003.1-200x also does not address the priorities of “system” processes. In general,
 5413 these processes should always execute in low-priority ranges to avoid conflict with other
 5414 realtime processes. Implementations should document the priority ranges in which system
 5415 processes run.

5416 The default scheduling policy is not defined. The effect of I/O interrupts and other system
 5417 processing activities is not defined. The temporary lending of priority from one process to
 5418 another (such as for the purposes of affecting freeing resources) by the system is not
 5419 addressed. Preemption of resources is not addressed. Restrictions on the ability of a
 5420 process to affect other processes beyond a certain level (influence levels) is not addressed.

5421 The rationale used to justify the simple time-quantum scheduler is that it is common
 5422 practice to depend upon this type of scheduling to ensure “fair” distribution of processor
 5423 resources among portions of the application that must interoperate in a serial fashion. Note
 5424 that IEEE Std 1003.1-200x is silent with respect to the setting of this time quantum, or
 5425 whether it is a system-wide value or a per-process value, although it appears that the
 5426 prevailing realtime practice is for it to be a system-wide value.

5427 In a system with N processes at a given priority, all processor-bound, in which the time
 5428 quantum is equal for all processes at a specific priority level, the following assumptions
 5429 are made of such a scheduling policy:

- 5430 1. A time quantum Q exists and the current process will own control of the processor
 5431 for at least a duration of Q and will have the processor for a duration of Q .
- 5432 2. The N th process at that priority will control a processor within a duration of $(N-1)$
 5433 $\times Q$.

5434 These assumptions are necessary to provide equal access to the processor and bounded
 5435 response from the application.

5436 The assumptions hold for the described scheduling policy only if no system overhead,
 5437 such as interrupt servicing, is present. If the interrupt servicing load is non-zero, then one
 5438 of the two assumptions becomes fallacious, based upon how Q is measured by the system.

5439 If Q is measured by clock time, then the assumption that the process obtains a duration Q
 5440 processor time is false if interrupt overhead exists. Indeed, a scenario can be constructed
 5441 with N processes in which a single process undergoes complete processor starvation if a
 5442 peripheral device, such as an analog-to-digital converter, generates significant interrupt

5443 activity periodically with a period of $N \times Q$.

5444 If Q is measured as actual processor time, then the assumption that the N th process runs in
5445 within the duration $(N-1) \times Q$ is false.

5446 It should be noted that SCHED_FIFO suffers from interrupt-based delay as well. However,
5447 for SCHED_FIFO, the implied response of the system is “as soon as possible”, so that the
5448 interrupt load for this case is a vendor selection and not a compliance issue.

5449 With this in mind, it is necessary either to complete the definition by including bounds on
5450 the interrupt load, or to modify the assumptions that can be made about the scheduling
5451 policy.

5452 Since the motivation of inclusion of the policy is common usage, and since current
5453 applications do not enjoy the luxury of bounded interrupt load, item (2) above is sufficient
5454 to express existing application needs and is less restrictive in the standard definition. No
5455 difference in interface is necessary.

5456 In an implementation in which the time quantum is equal for all processes at a specific
5457 priority, our assumptions can then be restated as:

- 5458 — A time quantum Q exists, and a processor-bound process will be rescheduled after a
5459 duration of, at most, Q . Time quantum Q may be defined in either wall clock time or
5460 execution time.
- 5461 — In general, the N th process of a priority level should wait no longer than $(N-1) \times Q$
5462 time to execute, assuming no processes exist at higher priority levels.
- 5463 — No process should wait indefinitely.

5464 For implementations supporting per-process time quanta, these assumptions can be
5465 readily extended.

5466 **Sporadic Server Scheduling Policy**

5467 The sporadic server is a mechanism defined for scheduling aperiodic activities in time-critical
5468 realtime systems. This mechanism reserves a certain bounded amount of execution capacity for
5469 processing aperiodic events at a high priority level. Any aperiodic events that cannot be
5470 processed within the bounded amount of execution capacity are executed in the background at a
5471 low priority level. Thus, a certain amount of execution capacity can be guaranteed to be
5472 available for processing periodic tasks, even under burst conditions in the arrival of aperiodic
5473 processing requests (that is, a large number of requests in a short time interval). The sporadic
5474 server also simplifies the schedulability analysis of the realtime system, because it allows
5475 aperiodic processes or threads to be treated as if they were periodic. The sporadic server was
5476 first described by Sprunt, et al.

5477 The key concept of the sporadic server is to provide and limit a certain amount of computation
5478 capacity for processing aperiodic events at their assigned normal priority, during a time interval
5479 called the “replenishment period”. Once the entity controlled by the sporadic server mechanism
5480 is initialized with its period and execution-time budget attributes, it preserves its execution
5481 capacity until an aperiodic request arrives. The request will be serviced (if there are no higher
5482 priority activities pending) as long as there is execution capacity left. If the request is completed,
5483 the actual execution time used to service it is subtracted from the capacity, and a replenishment
5484 of this amount of execution time is scheduled to happen one replenishment period after the
5485 arrival of the aperiodic request. If the request is not completed, because there is no execution
5486 capacity left, then the aperiodic process or thread is assigned a lower background priority. For
5487 each portion of consumed execution capacity the execution time used is replenished after one
5488 replenishment period. At the time of replenishment, if the sporadic server was executing at a

5489 background priority level, its priority is elevated to the normal level. Other similar
 5490 replenishment policies have been defined, but the one presented here represents a compromise
 5491 between efficiency and implementation complexity.

5492 The interface that appears in this section defines a new scheduling policy for threads and
 5493 processes that behaves according to the rules of the sporadic server mechanism. Scheduling
 5494 attributes are defined and functions are provided to allow the user to set and get the parameters
 5495 that control the scheduling behavior of this mechanism, namely the normal and low priority, the
 5496 replenishment period, the maximum number of pending replenishment operations, and the
 5497 initial execution-time budget.

5498 • Scheduling Aperiodic Activities

5499 Virtually all realtime applications are required to process aperiodic activities. In many
 5500 cases, there are tight timing constraints that the response to the aperiodic events must
 5501 meet. Usual timing requirements imposed on the response to these events are:

- 5502 — The effects of an aperiodic activity on the response time of lower priority activities
 5503 must be controllable and predictable.
- 5504 — The system must provide the fastest possible response time to aperiodic events.
- 5505 — It must be possible to take advantage of all the available processing bandwidth not
 5506 needed by time-critical activities to enhance average-case response times to aperiodic
 5507 events.

5508 Traditional methods for scheduling aperiodic activities are background processing, polling
 5509 tasks, and direct event execution:

- 5510 — Background processing consists of assigning a very low priority to the processing of
 5511 aperiodic events. It utilizes all the available bandwidth in the system that has not
 5512 been consumed by higher priority threads. However, it is very difficult, or
 5513 impossible, to meet requirements on average-case response time, because the
 5514 aperiodic entity has to wait for the execution of all other entities which have higher
 5515 priority.
- 5516 — Polling consists of creating a periodic process or thread for servicing aperiodic
 5517 requests. At regular intervals, the polling entity is started and its services
 5518 accumulated pending aperiodic requests. If no aperiodic requests are pending, the
 5519 polling entity suspends itself until its next period. Polling allows the aperiodic
 5520 requests to be processed at a higher priority level. However, worst and average-case
 5521 response times of polling entities are a direct function of the polling period, and there
 5522 is execution overhead for each polling period, even if no event has arrived. If the
 5523 deadline of the aperiodic activity is short compared to the inter-arrival time, the
 5524 polling frequency must be increased to guarantee meeting the deadline. For this case,
 5525 the increase in frequency can dramatically reduce the efficiency of the system and,
 5526 therefore, its capacity to meet all deadlines. Yet, polling represents a good way to
 5527 handle a large class of practical problems because it preserves system predictability,
 5528 and because the amortized overhead drops as load increases.
- 5529 — Direct event execution consists of executing the aperiodic events at a high fixed-
 5530 priority level. Typically, the aperiodic event is processed by an interrupt service
 5531 routine as soon as it arrives. This technique provides predictable response times for
 5532 aperiodic events, but makes the response times of all lower priority activities
 5533 completely unpredictable under burst arrival conditions. Therefore, if the density of
 5534 aperiodic event arrivals is unbounded, it may be a dangerous technique for time-
 5535 critical systems. Yet, for those cases in which the physics of the system imposes a
 5536 bound on the event arrival rate, it is probably the most efficient technique.

5537 — The sporadic server scheduling algorithm combines the predictability of the polling
 5538 approach with the short response times of the direct event execution. Thus, it allows
 5539 systems to meet an important class of application requirements that cannot be met by
 5540 using the traditional approaches. Multiple sporadic servers with different attributes
 5541 can be applied to the scheduling of multiple classes of aperiodic events, each with
 5542 different kinds of timing requirements, such as individual deadlines, average
 5543 response times, and so on. It also has many other interesting applications for
 5544 realtime, such as scheduling producer/consumer tasks in time-critical systems,
 5545 limiting the effects of faults on the estimation of task execution-time requirements,
 5546 and so on.

5547 • Existing Practice

5548 The sporadic server has been used in different kinds of applications, including military
 5549 avionics, robot control systems, industrial automation systems, and so on. There are
 5550 examples of many systems that cannot be successfully scheduled using the classic
 5551 approaches, such as direct event execution, or polling, and are schedulable using a
 5552 sporadic server scheduler. The sporadic server algorithm itself can successfully schedule
 5553 all systems scheduled with direct event execution or polling.

5554 The sporadic server scheduling policy has been implemented as a commercial product in
 5555 the run-time system of the Verdex Ada compiler. There are also many applications that
 5556 have used a much less efficient application-level sporadic server. These realtime
 5557 applications would benefit from a sporadic server scheduler implemented at the scheduler
 5558 level.

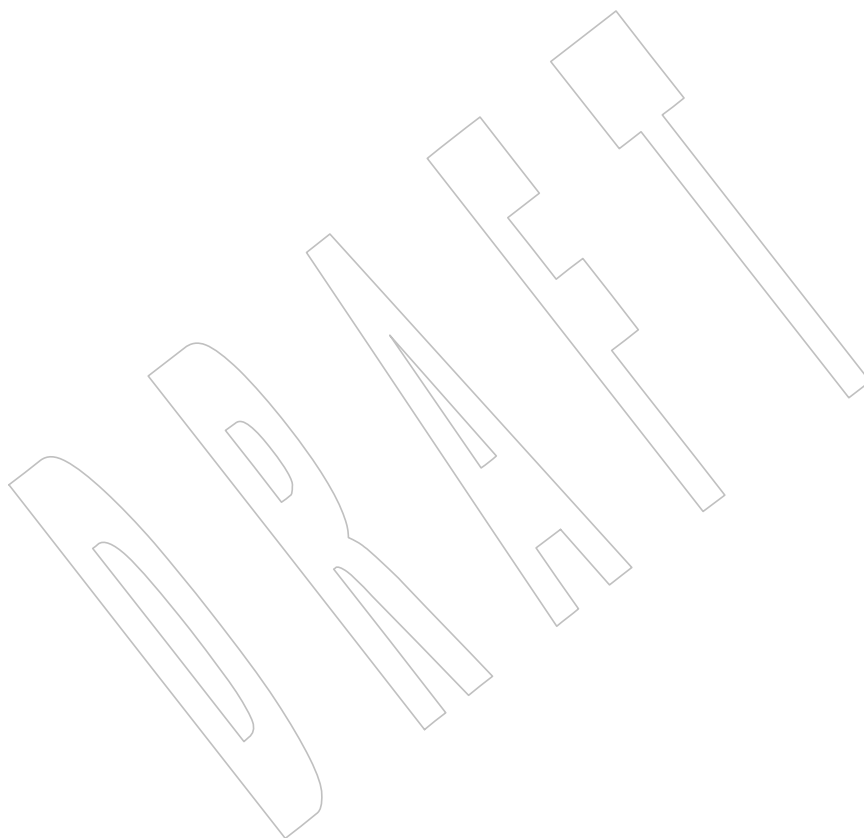
5559 • Library-Level *versus* Kernel-Level Implementation

5560 The sporadic server interface described in this section requires the sporadic server policy
 5561 to be implemented at the same level as the scheduler. This means that the process sporadic
 5562 server must be implemented at the kernel level and the thread sporadic server policy
 5563 implemented at the same level as the thread scheduler; that is, kernel or library level.

5564 In an earlier interface for the sporadic server, this mechanism was implementable at a
 5565 different level than the scheduler. This feature allowed the implementor to choose between
 5566 an efficient scheduler-level implementation, or a simpler user or library-level
 5567 implementation. However, the working group considered that this interface made the use
 5568 of sporadic servers more complex, and that library-level implementations would lack some
 5569 of the important functionality of the sporadic server, namely the limitation of the actual
 5570 execution time of aperiodic activities. The working group also felt that the interface
 5571 described in this chapter does not preclude library-level implementations of threads
 5572 intended to provide efficient low-overhead scheduling for those threads that are not
 5573 scheduled under the sporadic server policy.

5574 • Range of Scheduling Priorities

5575 Each of the scheduling policies supported in IEEE Std 1003.1-200x has an associated range
 5576 of priorities. The priority ranges for each policy might or might not overlap with the
 5577 priority ranges of other policies. For time-critical realtime applications it is usual for
 5578 periodic and aperiodic activities to be scheduled together in the same processor. Periodic
 5579 activities will usually be scheduled using the SCHED_FIFO scheduling policy, while
 5580 aperiodic activities may be scheduled using SCHED_SPORADIC. Since the application
 5581 developer will require complete control over the relative priorities of these activities in
 5582 order to meet his timing requirements, it would be desirable for the priority ranges of
 5583 SCHED_FIFO and SCHED_SPORADIC to overlap completely. Therefore, although
 5584 IEEE Std 1003.1-200x does not require any particular relationship between the different
 5585 priority ranges, it is recommended that these two ranges should coincide.



1. One-shot

A one-shot timer is a timer that is armed with an initial expiration time, either relative to the current time or at an absolute time (based on some timing base, such as time in seconds and nanoseconds since the Epoch). The timer expires once and then is disarmed. With the specified facilities, this is accomplished by setting the *it_value* member of the *value* argument to the desired expiration time and the *it_interval* member to zero.

2. Periodic

A periodic timer is a timer that is armed with an initial expiration time, again either relative or absolute, and a repetition interval. When the initial expiration occurs, the timer is reloaded with the repetition interval and continues counting. With the specified facilities, this is accomplished by setting the *it_value* member of the *value* argument to the desired initial expiration time and the *it_interval* member to the desired repetition interval.

For both of these types of timers, the time of the initial timer expiration can be specified in two ways:

1. Relative (to the current time)

2. Absolute

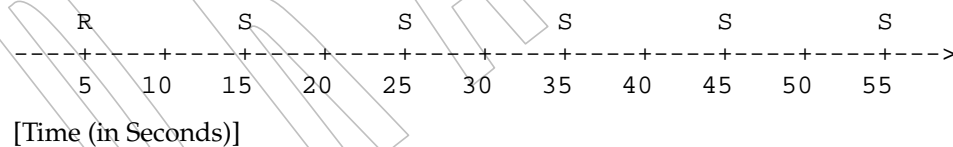
• Examples of Using Realtime Timers

In the diagrams below, *S* indicates a program schedule, *R* shows a schedule method request, and *E* suggests an internal operating system event.

— Periodic Timer: Data Logging

During an experiment, it might be necessary to log realtime data periodically to an internal buffer or to a mass storage device. With a periodic scheduling method, a logging module can be started automatically at fixed time intervals to log the data.

Program schedule is requested every 10 seconds.

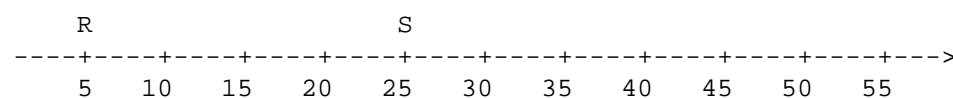


To achieve this type of scheduling using the specified facilities, one would allocate a per-process timer based on clock ID `CLOCK_REALTIME`. Then the timer would be armed via a call to `timer_settime()` with the `TIMER_ABSTIME` flag reset, and with an initial expiration value and a repetition interval of 10 seconds.

— One-shot Timer (Relative Time): Device Initialization

In an emission test environment, large sample bags are used to capture the exhaust from a vehicle. The exhaust is purged from these bags before each and every test. With a one-shot timer, a module could initiate the purge function and then suspend itself for a predetermined period of time while the sample bags are prepared.

Program schedule requested 20 seconds after call is issued.



5675 [Time (in Seconds)]

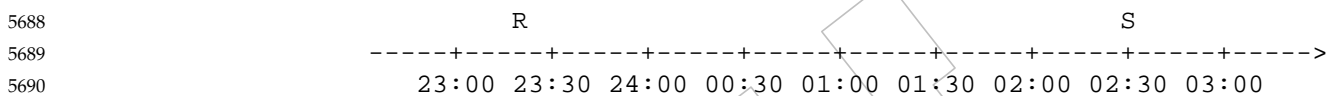
5676 To achieve this type of scheduling using the specified facilities, one would allocate a
5677 per-process timer based on clock ID `CLOCK_REALTIME`. Then the timer would be
5678 armed via a call to `timer_settime()` with the `TIMER_ABSTIME` flag reset, and with an
5679 initial expiration value of 20 seconds and a repetition interval of zero.

5680 Note that if the program wishes merely to suspend itself for the specified interval, it
5681 could more easily use `nanosleep()`.

5682 — One-shot Timer (Absolute Time): Data Transmission

5683 The results from an experiment are often moved to a different system within a
5684 network for postprocessing or archiving. With an absolute one-shot timer, a module
5685 that moves data from a test-cell computer to a host computer can be automatically
5686 scheduled on a daily basis.

5687 Program schedule requested for 2:30 a.m.



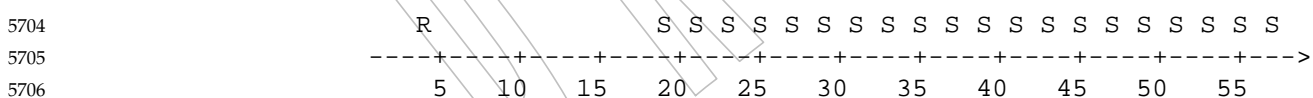
5691 [Time of Day]

5692 To achieve this type of scheduling using the specified facilities, a per-process timer
5693 would be allocated based on clock ID `CLOCK_REALTIME`. Then the timer would be
5694 armed via a call to `timer_settime()` with the `TIMER_ABSTIME` flag set, and an initial
5695 expiration value equal to 2:30 a.m. of the next day.

5696 — Periodic Timer (Relative Time): Signal Stabilization

5697 Some measurement devices, such as emission analyzers, do not respond
5698 instantaneously to an introduced sample. With a periodic timer with a relative initial
5699 expiration time, a module that introduces a sample and records the average response
5700 could suspend itself for a predetermined period of time while the signal is stabilized
5701 and then sample at a fixed rate.

5702 Program schedule requested 15 seconds after call is issued and every 2 seconds
5703 thereafter.



5707 [Time (in Seconds)]

5708 To achieve this type of scheduling using the specified facilities, one would allocate a
5709 per-process timer based on clock ID `CLOCK_REALTIME`. Then the timer would be
5710 armed via a call to `timer_settime()` with `TIMER_ABSTIME` flag reset, and with an
5711 initial expiration value of 15 seconds and a repetition interval of 2 seconds.

5712 — Periodic Timer (Absolute Time): Work Shift-related Processing

5713 Resource utilization data is useful when time to perform experiments is being
5714 scheduled at a facility. With a periodic timer with an absolute initial expiration time,
5715 a module can be scheduled at the beginning of a work shift to gather resource
5716 utilization data throughout the shift. This data can be used to allocate resources
5717 effectively to minimize bottlenecks and delays and maximize facility throughput.

5718 Program schedule requested for 2:00 a.m. and every 15 minutes thereafter.

```

5719             R                               S S S S S S
5720 -----+-----+-----+-----+-----+-----+-----+----->
5721             23:00 23:30 24:00 00:30 01:00 01:30 02:00 02:30 03:00

```

5722 [Time of Day]

5723 To achieve this type of scheduling using the specified facilities, one would allocate a
5724 per-process timer based on clock ID `CLOCK_REALTIME`. Then the timer would be
5725 armed via a call to `timer_settime()` with `TIMER_ABSTIME` flag set, and with an initial
5726 expiration value equal to 2:00 a.m. and a repetition interval equal to 15 minutes.

5727 • Relationship of Timers to Clocks

5728 The relationship between clocks and timers armed with an absolute time is
5729 straightforward: a timer expiration signal is requested when the associated clock reaches
5730 or exceeds the specified time. The relationship between clocks and timers armed with a
5731 relative time (an interval) is less obvious, but not unintuitive. In this case, a timer
5732 expiration signal is requested when the specified interval, *as measured by the associated clock*,
5733 has passed. For the required `CLOCK_REALTIME` clock, this allows timer expiration
5734 signals to be requested at specified “wall clock” times (absolute), or when a specified
5735 interval of “realtime” has passed (relative). For an implementation-defined clock—say, a
5736 process virtual time clock—timer expirations could be requested when the process has
5737 used a specified total amount of virtual time (absolute), or when it has used a specified
5738 *additional* amount of virtual time (relative).

5739 The interfaces also allow flexibility in the implementation of the functions. For example, an
5740 implementation could convert all absolute times to intervals by subtracting the clock value
5741 at the time of the call from the requested expiration time and “counting down” at the
5742 supported resolution. Or it could convert all relative times to absolute expiration time by
5743 adding in the clock value at the time of the call and comparing the clock value to the
5744 expiration time at the supported resolution. Or it might even choose to maintain absolute
5745 times as absolute and compare them to the clock value at the supported resolution for
5746 absolute timers, and maintain relative times as intervals and count them down at the
5747 resolution supported for relative timers. The choice will be driven by efficiency
5748 considerations and the underlying hardware or software clock implementation.

5749 • Data Definitions for Clocks and Timers

5750 IEEE Std 1003.1-200x uses a time representation capable of supporting nanosecond
5751 resolution timers for the following reasons:

- 5752 — To enable IEEE Std 1003.1-200x to represent those computer systems already using
5753 nanosecond or submicrosecond resolution clocks.
- 5754 — To accommodate those per-process timers that might need nanoseconds to specify an
5755 absolute value of system-wide clocks, even though the resolution of the per-process
5756 timer may only be milliseconds, or *vice versa*.
- 5757 — Because the number of nanoseconds in a second can be represented in 32 bits.

5758 Time values are represented in the `timespec` structure. The `tv_sec` member is of type `time_t`
5759 so that this member is compatible with time values used by POSIX.1 functions and the
5760 ISO C standard. The `tv_nsec` member is a **signed long** in order to simplify and clarify code
5761 that decrements or finds differences of time values. Note that because 1 billion (number of
5762 nanoseconds per second) is less than half of the value representable by a signed 32-bit
5763 value, it is always possible to add two valid fractional seconds represented as integral
5764 nanoseconds without overflowing the signed 32-bit value.

5765 A maximum allowable resolution for the `CLOCK_REALTIME` clock of 20 ms (1/50

5766 seconds) was chosen to allow line frequency clocks in European countries to be
 5767 conforming. 60 Hz clocks in the U.S. will also be conforming, as will finer granularity
 5768 clocks, although a Strictly Conforming Application cannot assume a granularity of less
 5769 than 20 ms (1/50 seconds).

5770 The minimum allowable maximum time allowed for the CLOCK_REALTIME clock and
 5771 the function *nanosleep()*, and timers created with *clock_id=CLOCK_REALTIME*, is
 5772 determined by the fact that the *tv_sec* member is of type **time_t**.

5773 IEEE Std 1003.1-200x specifies that timer expirations must not be delivered early, and
 5774 *nanosleep()* must not return early due to quantization error. IEEE Std 1003.1-200x discusses
 5775 the various implementations of *alarm()* in the rationale and states that implementations
 5776 that do not allow alarm signals to occur early are the most appropriate, but refrained from
 5777 mandating this behavior. Because of the importance of predictability to realtime
 5778 applications, IEEE Std 1003.1-200x takes a stronger stance.

5779 The developers of IEEE Std 1003.1-200x considered using a time representation that differs
 5780 from POSIX.1b in the second 32 bit of the 64-bit value. Whereas POSIX.1b defines this field
 5781 as a fractional second in nanoseconds, the other methodology defines this as a binary
 5782 fraction of one second, with the radix point assumed before the most significant bit.

5783 POSIX.1b is a software, source-level standard and most of the benefits of the alternate
 5784 representation are enjoyed by hardware implementations of clocks and algorithms. It was
 5785 felt that mandating this format for POSIX.1b clocks and timers would unnecessarily
 5786 burden the application writer with writing, possibly non-portable, multiple precision
 5787 arithmetic packages to perform conversion between binary fractions and integral units
 5788 such as nanoseconds, milliseconds, and so on.

5789 **Rationale for the Monotonic Clock**

5790 For those applications that use time services to achieve realtime behavior, changing the value of
 5791 the clock on which these services rely may cause erroneous timing behavior. For these
 5792 applications, it is necessary to have a monotonic clock which cannot run backwards, and which
 5793 has a maximum clock jump that is required to be documented by the implementation.
 5794 Additionally, it is desirable (but not required by IEEE Std 1003.1-200x) that the monotonic clock
 5795 increases its value uniformly. This clock should not be affected by changes to the system time;
 5796 for example, to synchronize the clock with an external source or to account for leap seconds.
 5797 Such changes would cause errors in the measurement of time intervals for those time services
 5798 that use the absolute value of the clock.

5799 One could argue that by defining the behavior of time services when the value of a clock is
 5800 changed, deterministic realtime behavior can be achieved. For example, one could specify that
 5801 relative time services should be unaffected by changes in the value of a clock. However, there are
 5802 time services that are based upon an absolute time, but that are essentially intended as relative
 5803 time services. For example, *pthread_cond_timedwait()* uses an absolute time to allow it to wake
 5804 up after the required interval despite spurious wakeups. Although sometimes the
 5805 *pthread_cond_timedwait()* timeouts are absolute in nature, there are many occasions in which they
 5806 are relative, and their absolute value is determined from the current time plus a relative time
 5807 interval. In this latter case, if the clock changes while the thread is waiting, the wait interval will
 5808 not be the expected length. If a *pthread_cond_timedwait()* function were created that would take a
 5809 relative time, it would not solve the problem because to retain the intended “deadline” a thread
 5810 would need to compensate for latency due to the spurious wakeup, and preemption between
 5811 wakeup and the next wait.

5812 The solution is to create a new monotonic clock, whose value does not change except for the
 5813 regular ticking of the clock, and use this clock for implementing the various relative timeouts

5814 that appear in the different POSIX interfaces, as well as allow *pthread_cond_timedwait()* to choose
 5815 this new clock for its timeout. A new *clock_nanosleep()* function is created to allow an application
 5816 to take advantage of this newly defined clock. Notice that the monotonic clock may be
 5817 implemented using the same hardware clock as the system clock.

5818 Relative timeouts for *sigtimedwait()* and *aio_suspend()* have been redefined to use the monotonic
 5819 clock, if present. The *alarm()* function has not been redefined, because the same effect but with
 5820 better resolution can be achieved by creating a timer (for which the appropriate clock may be
 5821 chosen).

5822 The *pthread_cond_timedwait()* function has been treated in a different way, compared to other
 5823 functions with absolute timeouts, because it is used to wait for an event, and thus it may have a
 5824 deadline, while the other timeouts are generally used as an error recovery mechanism, and for
 5825 them the use of the monotonic clock is not so important. Since the desired timeout for the
 5826 *pthread_cond_timedwait()* function may either be a relative interval or an absolute time of day
 5827 deadline, a new initialization attribute has been created for condition variables to specify the
 5828 clock that is used for measuring the timeout in a call to *pthread_cond_timedwait()*. In this way, if
 5829 a relative timeout is desired, the monotonic clock will be used; if an absolute deadline is
 5830 required instead, the *CLOCK_REALTIME* or another appropriate clock may be used. This
 5831 capability has not been added to other functions with absolute timeouts because for those
 5832 functions the expected use of the timeout is mostly to prevent errors, and not so often to meet
 5833 precise deadlines. As a consequence, the complexity of adding this capability is not justified by
 5834 its perceived application usage.

5835 The *nanosleep()* function has not been modified with the introduction of the monotonic clock.
 5836 Instead, a new *clock_nanosleep()* function has been created, in which the desired clock may be
 5837 specified in the function call.

5838 • History of Resolution Issues

5839 Due to the shift from relative to absolute timeouts in IEEE Std 1003.1d-1999, the
 5840 amendments to the *sem_timedwait()*, *pthread_mutex_timedlock()*, *mq_timedreceive()*, and
 5841 *mq_timedsend()* functions of that standard have been removed. Those amendments
 5842 specified that *CLOCK_MONOTONIC* would be used for the (relative) timeouts if the
 5843 Monotonic Clock option was supported.

5844 Having these functions continue to be tied solely to *CLOCK_MONOTONIC* would not
 5845 work. Since the absolute value of a time value obtained from *CLOCK_MONOTONIC* is
 5846 unspecified, under the absolute timeouts interface, applications would behave differently
 5847 depending on whether the Monotonic Clock option was supported or not (because the
 5848 absolute value of the clock would have different meanings in either case).

5849 Two options were considered:

- 5850 1. Leave the current behavior unchanged, which specifies the *CLOCK_REALTIME*
 5851 clock for these (absolute) timeouts, to allow portability of applications between
 5852 implementations supporting or not the Monotonic Clock option.
- 5853 2. Modify these functions in the way that *pthread_cond_timedwait()* was modified to
 5854 allow a choice of clock, so that an application could use *CLOCK_REALTIME* when
 5855 it is trying to achieve an absolute timeout and *CLOCK_MONOTONIC* when it is
 5856 trying to achieve a relative timeout.

5857 It was decided that the features of *CLOCK_MONOTONIC* are not as critical to these
 5858 functions as they are to *pthread_cond_timedwait()*. The *pthread_cond_timedwait()* function is
 5859 given a relative timeout; the timeout may represent a deadline for an event. When these
 5860 functions are given relative timeouts, the timeouts are typically for error recovery
 5861 purposes and need not be so precise.

5862 Therefore, it was decided that these functions should be tied to `CLOCK_REALTIME` and
5863 not complicated by being given a choice of clock.

5864 Execution Time Monitoring

5865 • Introduction

5866 The main goals of the execution time monitoring facilities defined in this chapter are to
5867 measure the execution time of processes and threads and to allow an application to
5868 establish CPU time limits for these entities.

5869 The analysis phase of time-critical realtime systems often relies on the measurement of
5870 execution times of individual threads or processes to determine whether the timing
5871 requirements will be met. Also, performance analysis techniques for soft deadline realtime
5872 systems rely heavily on the determination of these execution times. The execution time
5873 monitoring functions provide application developers with the ability to measure these
5874 execution times online and open the possibility of dynamic execution-time analysis and
5875 system reconfiguration, if required.

5876 The second goal of allowing an application to establish execution time limits for individual
5877 processes or threads and detecting when they overrun allows program robustness to be
5878 increased by enabling online checking of the execution times.

5879 If errors are detected—possibly because of erroneous program constructs, the existence of
5880 errors in the analysis phase, or a burst of event arrivals—online detection and recovery is
5881 possible in a portable way. This feature can be extremely important for many time-critical
5882 applications. Other applications require trapping CPU-time errors as a normal way to exit
5883 an algorithm; for instance, some realtime artificial intelligence applications trigger a
5884 number of independent inference processes of varying accuracy and speed, limit how long
5885 they can run, and pick the best answer available when time runs out. In many periodic
5886 systems, overrun processes are simply restarted in the next resource period, after necessary
5887 end-of-period actions have been taken. This allows algorithms that are inherently data-
5888 dependent to be made predictable.

5889 The interface that appears in this chapter defines a new type of clock, the CPU-time clock,
5890 which measures execution time. Each process or thread can invoke the clock and timer
5891 functions defined in POSIX.1 to use them. Functions are also provided to access the CPU-
5892 time clock of other processes or threads to enable remote monitoring of these clocks.
5893 Monitoring of threads of other processes is not supported, since these threads are not
5894 visible from outside of their own process with the interfaces defined in POSIX.1.

5895 • Execution Time Monitoring Interface

5896 The clock and timer interface defined in POSIX.1 historically only defined one clock, which
5897 measures wall-clock time. The requirements for measuring execution time of processes and
5898 threads, and setting limits to their execution time by detecting when they overrun, can be
5899 accomplished with that interface if a new kind of clock is defined. These new clocks
5900 measure execution time, and one is associated with each process and with each thread. The
5901 clock functions currently defined in POSIX.1 can be used to read and set these CPU-time
5902 clocks, and timers can be created using these clocks as their timing base. These timers can
5903 then be used to send a signal when some specified execution time has been exceeded. The
5904 CPU-time clocks of each process or thread can be accessed by using the symbols
5905 `CLOCK_PROCESS_CPUTIME_ID` or `CLOCK_THREAD_CPUTIME_ID`.

5906 The clock and timer interface defined in POSIX.1 and extended with the new kind of CPU-
5907 time clock would only allow processes or threads to access their own CPU-time clocks.
5908 However, many realtime systems require the possibility of monitoring the execution time

5909 of processes or threads from independent monitoring entities. In order to allow
 5910 applications to construct independent monitoring entities that do not require cooperation
 5911 from or modification of the monitored entities, two functions have been added:
 5912 *clock_getcpuclid()*, for accessing CPU-time clocks of other processes, and
 5913 *pthread_getcpuclid()*, for accessing CPU-time clocks of other threads. These functions
 5914 return the clock identifier associated with the process or thread specified in the call. These
 5915 clock IDs can then be used in the rest of the clock function calls.

5916 The clocks accessed through these functions could also be used as a timing base for the
 5917 creation of timers, thereby allowing independent monitoring entities to limit the CPU time
 5918 consumed by other entities. However, this possibility would imply additional complexity
 5919 and overhead because of the need to maintain a timer queue for each process or thread, to
 5920 store the different expiration times associated with timers created by different processes or
 5921 threads. The working group decided this additional overhead was not justified by
 5922 application requirements. Therefore, creation of timers attached to the CPU-time clocks of
 5923 other processes or threads has been specified as implementation-defined.

5924 • Overhead Considerations

5925 The measurement of execution time may introduce additional overhead in the thread
 5926 scheduling, because of the need to keep track of the time consumed by each of these
 5927 entities. In library-level implementations of threads, the efficiency of scheduling could be
 5928 somehow compromised because of the need to make a kernel call, at each context switch,
 5929 to read the process CPU-time clock. Consequently, a thread creation attribute called *cpu-*
 5930 *clock-requirement* was defined, to allow threads to disconnect their respective CPU-time
 5931 clocks. However, the Ballot Group considered that this attribute itself introduced some
 5932 overhead, and that in current implementations it was not worth the effort. Therefore, the
 5933 attribute was deleted, and thus thread CPU-time clocks are required for all threads if the
 5934 Thread CPU-Time Clocks option is supported.

5935 • Accuracy of CPU-Time Clocks

5936 The mechanism used to measure the execution time of processes and threads is specified in
 5937 IEEE Std 1003.1-200x as implementation-defined. The reason for this is that both the
 5938 underlying hardware and the implementation architecture have a very strong influence on
 5939 the accuracy achievable for measuring CPU time. For some implementations, the
 5940 specification of strict accuracy requirements would represent very large overheads, or even
 5941 the impossibility of being implemented.

5942 Since the mechanism for measuring execution time is implementation-defined, realtime
 5943 applications will be able to take advantage of accurate implementations using a portable
 5944 interface. Of course, strictly conforming applications cannot rely on any particular degree
 5945 of accuracy, in the same way as they cannot rely on a very accurate measurement of wall
 5946 clock time. There will always exist applications whose accuracy or efficiency requirements
 5947 on the implementation are more rigid than the values defined in IEEE Std 1003.1-200x or
 5948 any other standard.

5949 In any case, there is a minimum set of characteristics that realtime applications would
 5950 expect from most implementations. One such characteristic is that the sum of all the
 5951 execution times of all the threads in a process equals the process execution time, when no
 5952 CPU-time clocks are disabled. This need not always be the case because implementations
 5953 may differ in how they account for time during context switches. Another characteristic is
 5954 that the sum of the execution times of all processes in a system equals the number of
 5955 processors, multiplied by the elapsed time, assuming that no processor is idle during that
 5956 elapsed time. However, in some implementations it might not be possible to relate CPU
 5957 time to elapsed time. For example, in a heterogeneous multi-processor system in which

5958 each processor runs at a different speed, an implementation may choose to define each
5959 “second” of CPU time to be a certain number of “cycles” that a CPU has executed.

5960 • Existing Practice

5961 Measuring and limiting the execution time of each concurrent activity are common
5962 features of most industrial implementations of realtime systems. Almost all critical
5963 realtime systems are currently built upon a cyclic executive. With this approach, a regular
5964 timer interrupt kicks off the next sequence of computations. It also checks that the current
5965 sequence has completed. If it has not, then some error recovery action can be undertaken
5966 (or at least an overrun is avoided). Current software engineering principles and the
5967 increasing complexity of software are driving application developers to implement these
5968 systems on multi-threaded or multi-process operating systems. Therefore, if a POSIX
5969 operating system is to be used for this type of application, then it must offer the same level
5970 of protection.

5971 Execution time clocks are also common in most UNIX implementations, although these
5972 clocks usually have requirements different from those of realtime applications. The
5973 POSIX.1 *times()* function supports the measurement of the execution time of the calling
5974 process, and its terminated child processes. This execution time is measured in clock ticks
5975 and is supplied as two different values with the user and system execution times,
5976 respectively. BSD supports the function *getrusage()*, which allows the calling process to get
5977 information about the resources used by itself and/or all of its terminated child processes.
5978 The resource usage includes user and system CPU time. Some UNIX systems have options
5979 to specify high resolution (up to one microsecond) CPU-time clocks using the *times()* or
5980 the *getrusage()* functions.

5981 The *times()* and *getrusage()* interfaces do not meet important realtime requirements, such
5982 as the possibility of monitoring execution time from a different process or thread, or the
5983 possibility of detecting an execution time overrun. The latter requirement is supported in
5984 some UNIX implementations that are able to send a signal when the execution time of a
5985 process has exceeded some specified value. For example, BSD defines the functions
5986 *getitimer()* and *setitimer()*, which can operate either on a realtime clock (wall-clock), or on
5987 virtual-time or profile-time clocks which measure CPU time in two different ways. These
5988 functions do not support access to the execution time of other processes.

5989 IBM’s MVS operating system supports per-process and per-thread execution time clocks. It
5990 also supports limiting the execution time of a given process.

5991 Given all this existing practice, the working group considered that the POSIX.1 clocks and
5992 timers interface was appropriate to meet most of the requirements that realtime
5993 applications have for execution time clocks. Functions were added to get the CPU time
5994 clock IDs, and to allow/disallow the thread CPU-time clocks (in order to preserve the
5995 efficiency of some implementations of threads).

5996 • Clock Constants

5997 The definition of the manifest constants `CLOCK_PROCESS_CPUTIME_ID` and
5998 `CLOCK_THREAD_CPUTIME_ID` allows processes or threads, respectively, to access their
5999 own execution-time clocks. However, given a process or thread, access to its own
6000 execution-time clock is also possible if the clock ID of this clock is obtained through a call
6001 to *clock_getcpuclockid()* or *pthread_getcpuclockid()*. Therefore, these constants are not
6002 necessary and could be deleted to make the interface simpler. Their existence saves one
6003 system call in the first access to the CPU-time clock of each process or thread. The working
6004 group considered this issue and decided to leave the constants in IEEE Std 1003.1-200x
6005 because they are closer to the POSIX.1b use of clock identifiers.

6006
6007
6008
6009
6010
6011
6012
6013
6014
6015
6016
6017
6018
6019
6020
6021
6022
6023
6024
6025
6026
6027
6028
6029
6030
6031
6032
6033
6034
6035
6036
6037
6038
6039
6040
6041
6042
6043
6044
6045
6046
6047
6048
6049
6050
6051
6052
6053

- Library Implementations of Threads

In library implementations of threads, kernel entities and library threads can coexist. In this case, if the CPU-time clocks are supported, most of the clock and timer functions will need to have two implementations: one in the thread library, and one in the system calls library. The main difference between these two implementations is that the thread library implementation will have to deal with clocks and timers that reside in the thread space, while the kernel implementation will operate on timers and clocks that reside in kernel space. In the library implementation, if the clock ID refers to a clock that resides in the kernel, a kernel call will have to be made. The correct version of the function can be chosen by specifying the appropriate order for the libraries during the link process.

- History of Resolution Issues: Deletion of the *enable* Attribute

In early proposals, consideration was given to inclusion of an attribute called *enable* for CPU-time clocks. This would allow implementations to avoid the overhead of measuring execution time for those processes or threads for which this measurement was not required. However, this is unnecessary since processes are already required to measure execution time by the POSIX.1 *times()* function. Consequently, the *enable* attribute is not present.

Rationale Relating to Timeouts

- Requirements for Timeouts

Realtime systems which must operate reliably over extended periods without human intervention are characteristic in embedded applications such as avionics, machine control, and space exploration, as well as more mundane applications such as cable TV, security systems, and plant automation. A multi-tasking paradigm, in which many independent and/or cooperating software functions relinquish the processor(s) while waiting for a specific stimulus, resource, condition, or operation completion, is very useful in producing well engineered programs for such systems. For such systems to be robust and fault-tolerant, expected occurrences that are unduly delayed or that never occur must be detected so that appropriate recovery actions may be taken. This is difficult if there is no way for a task to regain control of a processor once it has relinquished control (blocked) awaiting an occurrence which, perhaps because of corrupted code, hardware malfunction, or latent software bugs, will not happen when expected. Therefore, the common practice in realtime operating systems is to provide a capability to time out such blocking services. Although there are several methods to achieve this already defined by POSIX, none are as reliable or efficient as initiating a timeout simultaneously with initiating a blocking service. This is especially critical in hard-realtime embedded systems because the processors typically have little time reserve, and allowed fault recovery times are measured in milliseconds rather than seconds.

The working group largely agreed that such timeouts were necessary and ought to become part of IEEE Std 1003.1-200x, particularly vendors of realtime operating systems whose customers had already expressed a strong need for timeouts. There was some resistance to inclusion of timeouts in IEEE Std 1003.1-200x because the desired effect, fault tolerance, could, in theory, be achieved using existing facilities and alternative software designs, but there was no compelling evidence that realtime system designers would embrace such designs at the sacrifice of performance and/or simplicity.

- Which Services should be Timed Out?

Originally, the working group considered the prospect of providing timeouts on all blocking services, including those currently existing in POSIX.1, POSIX.1b, and POSIX.1c, and future interfaces to be defined by other working groups, as sort of a general policy.

6054 This was rather quickly rejected because of the scope of such a change, and the fact that
 6055 many of those services would not normally be used in a realtime context. More traditional
 6056 timesharing solutions to timeout would suffice for most of the POSIX.1 interfaces, while
 6057 others had asynchronous alternatives which, while more complex to utilize, would be
 6058 adequate for some realtime and all non-realtime applications.

6059 The list of potential candidates for timeouts was narrowed to the following for further
 6060 consideration:

- 6061 — POSIX.1b
 - 6062 — *sem_wait()*
 - 6063 — *mq_receive()*
 - 6064 — *mq_send()*
 - 6065 — *lio_listio()*
 - 6066 — *aio_suspend()*
 - 6067 — *sigwait()* (timeout already implemented by *sigtimedwait()*)
- 6068 — POSIX.1c
 - 6069 — *pthread_mutex_lock()*
 - 6070 — *pthread_join()*
 - 6071 — *pthread_cond_wait()*
 - 6072 (timeout already implemented by *pthread_cond_timedwait()*)
- 6073 — POSIX.1
 - 6074 — *read()*
 - 6075 — *write()*

6076 After further review by the working group, the *lio_listio()*, *read()*, and *write()* functions (all
 6077 forms of blocking synchronous I/O) were eliminated from the list because of the
 6078 following:

- 6079 — Asynchronous alternatives exist
- 6080 — Timeouts can be implemented, albeit non-portably, in device drivers
- 6081 — A strong desire not to introduce modifications to POSIX.1 interfaces

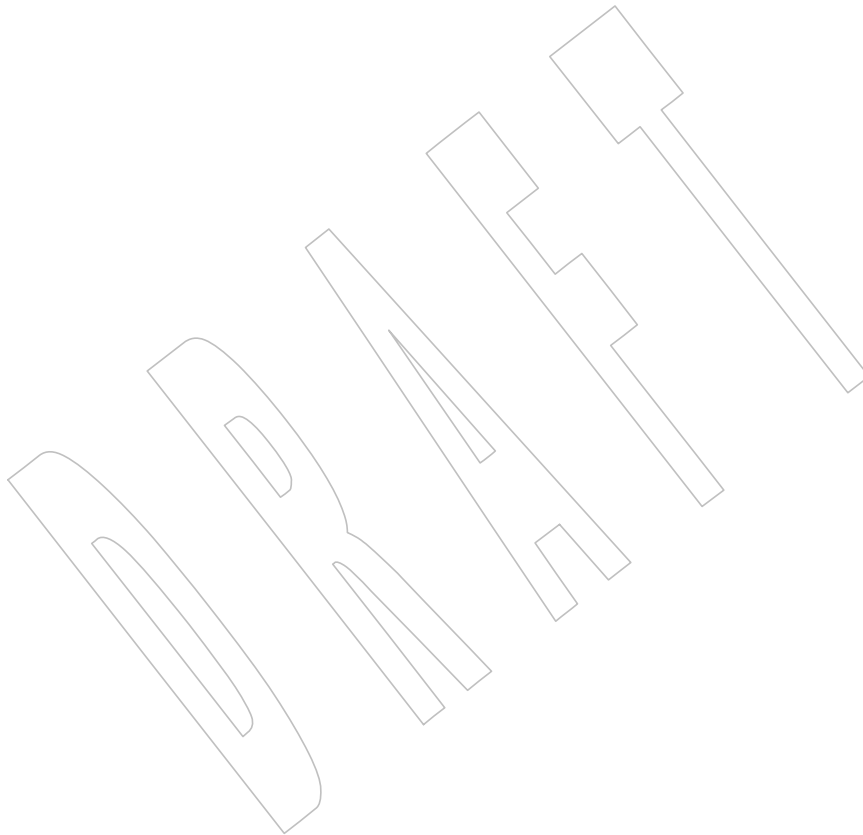
6082 The working group ultimately rejected *pthread_join()* since both that interface and a timed
 6083 variant of that interface are non-minimal and may be implemented as a function. See
 6084 below for a library implementation of *pthread_join()*.

6085 Thus, there was a consensus among the working group members to add timeouts to 4 of
 6086 the remaining 5 functions (the timeout for *aio_suspend()* was ultimately added directly to
 6087 POSIX.1b, while the others were added by POSIX.1d). However, *pthread_mutex_lock()*
 6088 remained contentious.

6089 Many feel that *pthread_mutex_lock()* falls into the same class as the other functions; that is,
 6090 it is desirable to time out a mutex lock because a mutex may fail to be unlocked due to
 6091 errant or corrupted code in a critical section (looping or branching outside of the unlock
 6092 code), and therefore is equally in need of a reliable, simple, and efficient timeout. In fact,
 6093 since mutexes are intended to guard small critical sections, most *pthread_mutex_lock()* calls
 6094 would be expected to obtain the lock without blocking nor utilizing any kernel service,
 6095 even in implementations of threads with global contention scope; the timeout alternative

need only be considered after it is determined that the thread must block.

Those opposed to timing out mutexes feel that the very simplicity of the mutex is compromised by adding a timeout semantic, and that to do so is senseless. They claim that if a timed mutex is really deemed useful by a particular application, then it can be constructed from the facilities already in POSIX.1b and POSIX.1c. The following two C-language library implementations of mutex locking with timeout represent the solutions offered (in both implementations, the timeout parameter is specified as absolute time, not



```

6144     {
6145     int timedout=FALSE;
6146     int error_status;
6147
6148     pthread_mutex_lock(&tm->mutex);
6149
6150     while (tm->locked && !timedout)
6151     {
6152         if ((error_status=pthread_cond_timedwait(&tm->cond,
6153         &tm->mutex,
6154         timeout))!=0)
6155         {
6156             if (error_status==ETIMEDOUT) timedout = TRUE;
6157         }
6158     }
6159
6160     if(timedout)
6161     {
6162         pthread_mutex_unlock(&tm->mutex);
6163         return ETIMEDOUT;
6164     }
6165     else
6166     {
6167         tm->locked = TRUE;
6168         pthread_mutex_unlock(&tm->mutex);
6169         return 0;
6170     }
6171 }
6172
6173 void timed_mutex_unlock(timed_mutex_t *tm)
6174 {
6175     pthread_mutex_lock(&tm->mutex); / for case assignment not atomic /
6176     tm->locked = FALSE;
6177     pthread_mutex_unlock(&tm->mutex);
6178     pthread_cond_signal(&tm->cond);
6179 }

```

The Condition Wait implementation effectively substitutes the *pthread_cond_timedwait()* function (which is currently timed out) for the desired *pthread_mutex_timedlock()*. Since waits on condition variables currently do not include protocols which avoid priority inversion, this method is generally unsuitable for realtime applications because it does not provide the same priority inversion protection as the untimed *pthread_mutex_lock()*. Also, for any given implementations of the current mutex and condition variable primitives, this library implementation has a performance cost at least 2.5 times that of the untimed *pthread_mutex_lock()* even in the case where the timed mutex is readily locked without blocking (the interfaces required for this case are shown in bold). Even in uniprocessors or where assignment is atomic, at least an additional *pthread_cond_signal()* is required. *pthread_mutex_timedlock()* could be implemented at effectively no performance penalty in this case because the timeout parameters need only be considered after it is determined that the mutex cannot be locked immediately.

Thus it has not yet been shown that the full semantics of mutex locking with timeout can be efficiently and reliably achieved using existing interfaces. Even if the existence of an acceptable library implementation were proven, it is difficult to justify why the interface itself should not be made portable, especially considering approval for the other four timeouts.

- Rationale for Library Implementation of *pthread_timedjoin()*

Library implementation of *pthread_timedjoin()*:

```

6194
6195
6196 /*
6197  * Construct a thread variety entirely from existing functions
6198  * with which a join can be done, allowing the join to time out.
6199  */
6200 #include <pthread.h>
6201 #include <time.h>
6202
6203 struct timed_thread {
6204     pthread_t t;
6205     pthread_mutex_t m;
6206     int exiting;
6207     pthread_cond_t exit_c;
6208     void *(*start_routine)(void *arg);
6209     void *arg;
6210     void *status;
6211 };
6212
6213 typedef struct timed_thread *timed_thread_t;
6214 static pthread_key_t timed_thread_key;
6215 static pthread_once_t timed_thread_once = PTHREAD_ONCE_INIT;
6216
6217 static void timed_thread_init()
6218 {
6219     pthread_key_create(&timed_thread_key, NULL);
6220 }
6221
6222 static void *timed_thread_start_routine(void *args)
6223 {
6224     /*
6225      * Routine to establish thread-specific data value and run the actual
6226      * thread start routine which was supplied to timed_thread_create().
6227      */
6228     {
6229         timed_thread_t tt = (timed_thread_t) args;
6230
6231         pthread_once(&timed_thread_once, timed_thread_init);
6232         pthread_setspecific(timed_thread_key, (void *)tt);
6233         timed_thread_exit((tt->start_routine)(tt->arg));
6234     }
6235
6236 int timed_thread_create(timed_thread_t ttp, const pthread_attr_t *attr,
6237 void *(*start_routine)(void *), void *arg)
6238
6239 /*
6240  * Allocate a thread which can be used with timed_thread_join().
6241  */
6242 {
6243     timed_thread_t tt;
6244     int result;
6245
6246     tt = (timed_thread_t) malloc(sizeof(struct timed_thread));
6247     pthread_mutex_init(&tt->m, NULL);
6248     tt->exiting = FALSE;
6249     pthread_cond_init(&tt->exit_c, NULL);

```

```

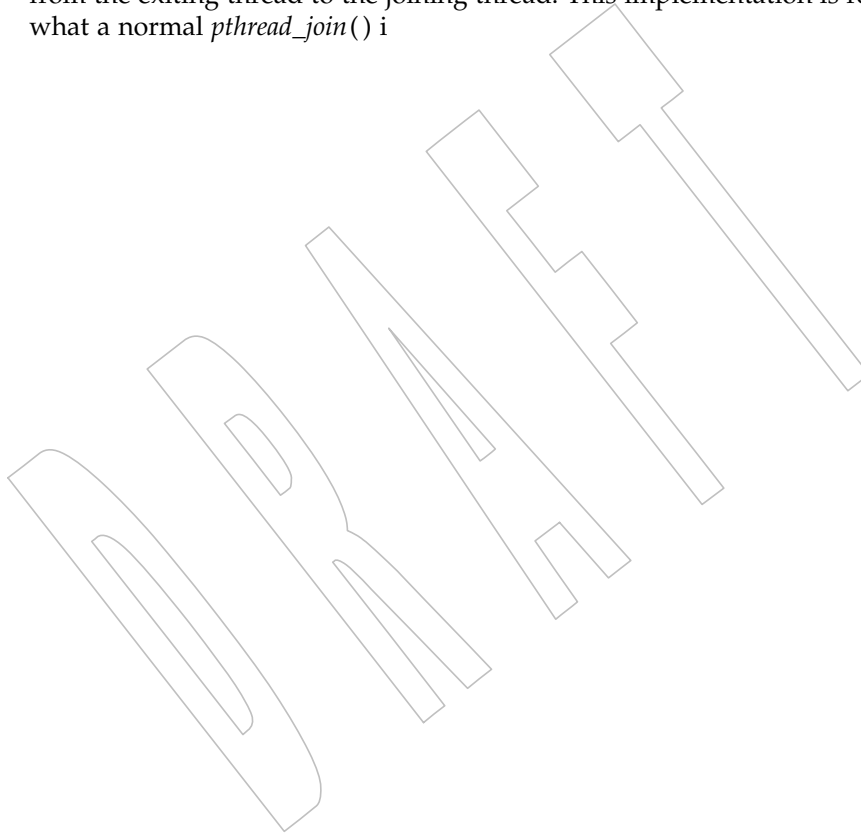
6241         tt->start_routine = start_routine;
6242         tt->arg = arg;
6243         tt->status = NULL;
6244
6245         if ((result = pthread_create(&tt->t, attr,
6246             timed_thread_start_routine, (void *)tt)) != 0) {
6247             free(tt);
6248             return result;
6249         }
6250
6251         pthread_detach(tt->t);
6252         ttp = tt;
6253         return 0;
6254     }
6255
6256     int timed_thread_join(timed_thread_t tt,
6257         struct timespec *timeout,
6258         void **status)
6259     {
6260         int result;
6261
6262         pthread_mutex_lock(&tt->m);
6263         result = 0;
6264         /*
6265          * Wait until the thread announces that it is exiting,
6266          * or until timeout.
6267          */
6268         while (result == 0 && ! tt->exiting) {
6269             result = pthread_cond_timedwait(&tt->exit_c, &tt->m, timeout);
6270         }
6271         pthread_mutex_unlock(&tt->m);
6272         if (result == 0 && tt->exiting) {
6273             *status = tt->status;
6274             free((void *)tt);
6275             return result;
6276         }
6277         return result;
6278     }
6279
6280     void timed_thread_exit(void *status)
6281     {
6282         timed_thread_t tt;
6283         void *specific;
6284
6285         if ((specific=pthread_getspecific(timed_thread_key)) == NULL){
6286             /*
6287              * Handle cases which won't happen with correct usage.
6288              */
6289             pthread_exit( NULL);
6290         }
6291         tt = (timed_thread_t) specific;
6292         pthread_mutex_lock(&tt->m);
6293         /*
6294          * Tell a joiner that we're exiting.
6295          */
6296         tt->status = status;

```



```
tt->exiting = TRUE;
pthread_cond_signal(&tt->exit_c);
pthread_mutex_unlock(&tt->m);
/*
 * Call pthread_exit() to call destructors and really
 * exit the thread.
 */
pthread_exit(NULL);
}
```

The `pthread_join()` C-language example shown above demonstrates that it is possible, using existing pthread facilities, to construct a variety of thread which allows for joining such a thread, but which allows the join operation to time out. It does this by using a `pthread_cond_timedwait()` to wait for the thread to exit. A **timed_thread_t** descriptor structure is used to pass parameters from the creating thread to the created thread, and from the exiting thread to the joining thread. This implementation is roughly equivalent to what a normal `pthread_join()` i



6340 that specifies absolute timeouts. In this case, the clock would have to be read to calculate
6341 the absolute wake-up time as the sum of the current time plus the relative timeout interval.
6342 In this case, if the thread is preempted after reading the clock but before making the timed-
6343 out call, the thread would be awakened earlier than desired.

6344 But the race condition with the absolute timeouts interface is not as bad as the one that
6345 happens with the relative timeout interface, because there are simple workarounds. For the
6346 absolute timeouts interface, if the timing requirement is a deadline, the deadline can still
6347 be met because the thread woke up earlier than the deadline. If the timeout is just used as
6348 an error recovery mechanism, the precision of timing is not really important. If the timing
6349 requirement is that between actions A and B a minimum interval of time must elapse, the
6350 absolute timeout interface can be safely used by reading the clock after action A has been
6351 started. It could be argued that, since the call with the absolute timeout is atomic from the
6352 application point of view, it is not possible to read the clock after action A, if this action is
6353 part of the timed-out call. But looking at the nature of the calls for which timeouts are
6354 specified (locking a mutex, waiting for a semaphore, waiting for a message, or waiting
6355 until there is space in a message queue), the timeouts that an application would build on
6356 these actions would not be triggered by these actions themselves, but by some other
6357 external action. For example, if waiting for a message to arrive to a message queue, and
6358 waiting for at least 20 milliseconds, this time interval would start to be counted from some
6359 event that would trigger both the action that produces the message, as well as the action
6360 that waits for the message to arrive, and not by the wait-for-message operation itself. In
6361 this case, the workaround proposed above could be used.

6362 For these reasons, the absolute timeout is preferred over the relative timeout interface.

6363 B.2.9 Threads

6364 Threads will normally be more expensive than subroutines (or functions, routines, and so on) if
6365 specialized hardware support is not provided. Nevertheless, threads should be sufficiently
6366 efficient to encourage their use as a medium to fine-grained structuring mechanism for
6367 parallelism in an application. Structuring an application using threads then allows it to take
6368 immediate advantage of any underlying parallelism available in the host environment. This
6369 means implementors are encouraged to optimize for fast execution at the possible expense of
6370 efficient utilization of storage. For example, a common thread creation technique is to cache
6371 appropriate thread data structures. That is, rather than releasing system resources, the
6372 implementation retains these resources and reuses them when the program next asks to create a
6373 new thread. If this reuse of thread resources is to be possible, there has to be very little unique
6374 state associated with each thread, because any such state has to be reset when the thread is
6375 reused.

6376 Thread Creation Attributes

6377 Attributes objects are provided for threads, mutexes, and condition variables as a mechanism to
6378 support probable future standardization in these areas without requiring that the interface itself
6379 be changed.

6380 Attributes objects provide clean isolation of the configurable aspects of threads. For example,
6381 “stack size” is an important attribute of a thread, but it cannot be expressed portably. When
6382 porting a threaded program, stack sizes often need to be adjusted. The use of attributes objects
6383 can help by allowing the changes to be isolated in a single place, rather than being spread across
6384 every instance of thread creation.

6385 Attributes objects can be used to set up *classes* of threads with similar attributes; for example,
6386 “threads with large stacks and high priority” or “threads with minimal stacks”. These classes
6387 can be defined in a single place and then referenced wherever threads need to be created.

6388 Changes to “class” decisions become straightforward, and detailed analysis of each
6389 *pthread_create()* call is not required.

6390 The attributes objects are defined as opaque types as an aid to extensibility. If these objects had
6391 been specified as structures, adding new attributes would force recompilation of all multi-
6392 threaded programs when the attributes objects are extended; this might not be possible if
6393 different program components were supplied by different vendors.

6394 Additionally, opaque attributes objects present opportunities for improving performance.
6395 Argument validity can be checked once when attributes are set, rather than each time a thread is
6396 created. Implementations will often need to cache kernel objects that are expensive to create.
6397 Opaque attributes objects provide an efficient mechanism to detect when cached objects become
6398 invalid due to attribute changes.

6399 Because assignment is not necessarily defined on a given opaque type, implementation-defined
6400 default values cannot be defined in a portable way. The solution to this problem is to allow
6401 attribute objects to be initialized dynamically by attributes object initialization functions, so that
6402 default values can be supplied automatically by the implementation.

6403 The following proposal was provided as a suggested alternative to the supplied attributes:

- 6404 1. Maintain the style of passing a parameter formed by the bitwise-inclusive OR of flags to
6405 the initialization routines (*pthread_create()*, *pthread_mutex_init()*, *pthread_cond_init()*). The
6406 parameter containing the flags should be an opaque type for extensibility. If no flags are
6407 set in the parameter, then the objects are created with default characteristics. An
6408 implementation may specify implementation-defined flag values and associated
6409 behavior.
- 6410 2. If further specialization of mutexes and condition variables is necessary, implementations
6411 may specify additional procedures that operate on the **pthread_mutex_t** and
6412 **pthread_cond_t** objects (instead of on attributes objects).

6413 The difficulties with this solution are:

- 6414 1. A bitmask is not opaque if bits have to be set into bit-vector attributes objects using
6415 explicitly-coded bitwise-inclusive OR operations. If the set of options exceeds an **int**,
6416 application programmers need to know the location of each bit. If bits are set or read by
6417 encapsulation (that is, *get*()* or *set*()* functions), then the bitmask is merely an
6418 implementation of attributes objects as currently defined and should not be exposed to
6419 the programmer.
- 6420 2. Many attributes are not Boolean or very small integral values. For example, scheduling
6421 policy may be placed in 3 bits or 4 bits, but priority requires 5 bits or more, thereby taking
6422 up at least 8 bits out of a possible 16 bits on machines with 16-bit integers. Because of this,
6423 the bitmask can only reasonably control whether particular attributes are set or not, and it
6424 cannot serve as the repository of the value itself. The value needs to be specified as a
6425 function parameter (which is non-extensible), or by setting a structure field (which is non-
6426 opaque), or by *get*()* and *set*()* functions (making the bitmask a redundant addition to
6427 the attributes objects).

6428 Stack size is defined as an optional attribute because the very notion of a stack is inherently
6429 machine-dependent. Some implementations may not be able to change the size of the stack, for
6430 example, and others may not need to because stack pages may be discontinuous and can be
6431 allocated and released on demand.

6432 The attribute mechanism has been designed in large measure for extensibility. Future extensions
6433 to the attribute mechanism or to any attributes object defined in IEEE Std 1003.1-200x have to be
6434 done with care so as not to affect binary-compatibility.

6435 Attribute objects, even if allocated by means of dynamic allocation functions such as *malloc()*,
6436 may have their size fixed at compile time. This means, for example, a *pthread_create()* in an
6437 implementation with extensions to the **pthread_attr_t** cannot look beyond the area that the
6438 binary application assumes is valid. This suggests that implementations should maintain a size
6439 field in the attributes object, as well as possibly version information, if extensions in different
6440 directions (possibly by different vendors) are to be accommodated.

6441 Thread Implementation Models

6442 There are various thread implementation models. At one end of the spectrum is the “library-
6443 thread model”. In such a model, the threads of a process are not visible to the operating system
6444 kernel, and the threads are not kernel-scheduled entities. The process is the only kernel-
6445 scheduled entity. The process is scheduled onto the processor by the kernel according to the
6446 scheduling attributes of the process. The threads are scheduled onto the single kernel-scheduled
6447 entity (the process) by the runtime library according to the scheduling attributes of the threads.
6448 A problem with this model is that it constrains concurrency. Since there is only one kernel-
6449 scheduled entity (namely, the process), only one thread per process can execute at a time. If the
6450 thread that is executing blocks on I/O, then the whole process blocks.

6451 At the other end of the spectrum is the “kernel-thread model”. In this model, all threads are
6452 visible to the operating system kernel. Thus, all threads are kernel-scheduled entities, and all
6453 threads can concurrently execute. The threads are scheduled onto processors by the kernel
6454 according to the scheduling attributes of the threads. The drawback to this model is that the
6455 creation and management of the threads entails operating system calls, as opposed to subroutine
6456 calls, which makes kernel threads heavier weight than library threads.

6457 Hybrids of these two models are common. A hybrid model offers the speed of library threads
6458 and the concurrency of kernel threads. In hybrid models, a process has some (relatively small)
6459 number of kernel scheduled entities associated with it. It also has a potentially much larger
6460 number of library threads associated with it. Some library threads may be bound to kernel-
6461 scheduled entities, while the other library threads are multiplexed onto the remaining kernel-
6462 scheduled entities. There are two levels of thread scheduling:

- 6463 1. The runtime library manages the scheduling of (unbound) library threads onto kernel-
6464 scheduled entities.
- 6465 2. The kernel manages the scheduling of kernel-scheduled entities onto processors.

6466 For this reason, a hybrid model is referred to as a two-level threads scheduling model. In this
6467 model, the process can have multiple concurrently executing threads; specifically, it can have as
6468 many concurrently executing threads as it has kernel-scheduled entities.

6469 Thread-Specific Data

6470 Many applications require that a certain amount of context be maintained on a per-thread basis
6471 across procedure calls. A common example is a multi-threaded library routine that allocates
6472 resources from a common pool and maintains an active resource list for each thread. The thread-
6473 specific data interface provided to meet these needs may be viewed as a two-dimensional array
6474 of values with keys serving as the row index and thread IDs as the column index (although the
6475 implementation need not work this way).

6476
6477
6478
6479
6480
6481
6482
6483
6484
6485
6486
6487
6488
6489
6490
6491
6492
6493
6494
6495
6496
6497
6498
6499
6500
6501
6502
6503
6504
6505
6506
6507
6508
6509
6510
6511
6512
6513
6514
6515
6516
6517
6518
6519
6520
6521

- Models

Three possible thread-specific data models were considered:

1. No Explicit Support

A standard thread-specific data interface is not strictly necessary to support applications that require per-thread context. One could, for example, provide a hash function that converted a `pthread_t` into an integer value that could then be used to index into a global array of per-thread data pointers. This hash function, in conjunction with `pthread_self()`, would be all the interface required to support a mechanism of this sort. Unfortunately, this technique is cumbersome. It can lead to duplicated code as each set of cooperating modules implements their own per-thread data management schemes.

2. Single (`void *`) Pointer

Another technique would be to provide a single word of per-thread storage and a pair of functions to fetch and store the value of this word. The word could then hold a pointer to a block of per-thread memory. The allocation, partitioning, and general use of this memory would be entirely up to the application. Although this method is not as problematic as technique 1, it suffers from interoperability problems. For example, all modules using the per-thread pointer would have to agree on a common usage protocol.

3. Key/Value Mechanism

This method associates an opaque key (for example, stored in a variable of type `pthread_key_t`) with each per-thread datum. These keys play the role of identifiers for per-thread data. This technique is the most generic and avoids the problems noted above, albeit at the cost of some complexity.

The primary advantage of the third model is its information hiding properties. Modules using this model are free to create and use their own key(s) independent of all other such usage, whereas the other models require that all modules that use thread-specific context explicitly cooperate with all other such modules. The data-independence provided by the third model is worth the additional interface.

- Requirements

It is important that it be possible to implement the thread-specific data interface without the use of thread private memory. To do otherwise would increase the weight of each thread, thereby limiting the range of applications for which the threads interfaces provided by IEEE Std 1003.1-200x is appropriate.

The values that one binds to the key via `pthread_setspecific()` may, in fact, be pointers to shared storage locations available to all threads. It is only the key/value bindings that are maintained on a per-thread basis, and these can be kept in any portion of the address space that is reserved for use by the calling thread (for example, on the stack). Thus, no per-thread MMU state is required to implement the interface. On the other hand, there is nothing in the interface specification to preclude the use of a per-thread MMU state if it is available (for example, the key values returned by `pthread_key_create()` could be thread private memory addresses).

- Standardization Issues

Thread-specific data is a requirement for a usable thread interface. The binding described in this section provides a portable thread-specific data mechanism for languages that do not directly support a thread-specific storage class. A binding to IEEE Std 1003.1-200x for a

6522 language that does include such a storage class need not provide this specific interface.

6523 If a language were to include the notion of thread-specific storage, it would be desirable
 6524 (but *not* required) to provide an implementation of the pthreads thread-specific data
 6525 interface based on the language feature. For example, assume that a compiler for a C-like
 6526 language supports a *private* storage class that provides thread-specific storage. Something
 6527 similar to the following macros might be used to effect a compatible implementation:

```
6528 #define pthread_key_t                private void *
6529 #define pthread_key_create(key)      /* no-op */
6530 #define pthread_setspecific(key,value) (key)=(value)
6531 #define pthread_getspecific(key)    (key)
```

6532 **Note:** For the sake of clarity, this example ignores destructor functions. A correct
 6533 implementation would have to support them.

6534 Barriers

- 6535 • Background

6536 Barriers are typically used in parallel DO/FOR loops to ensure that all threads have
 6537 reached a particular stage in a parallel computation before allowing any to proceed to the
 6538 next stage. Highly efficient implementation is possible on machines which support a
 6539 “Fetch and Add” operation as described in the referenced Almasi and Gottlieb (1989).

6540 The use of return value PTHREAD_BARRIER_SERIAL_THREAD is shown in the
 6541 following example:

```
6542 if ( (status=pthread_barrier_wait(&barrier)) ==
6543     PTHREAD_BARRIER_SERIAL_THREAD) {
6544     ...serial section
6545 }
6546     else if (status != 0) {
6547     ...error processing
6548 }
6549 status=pthread_barrier_wait(&barrier);
6550 ...
```

6551 This behavior allows a serial section of code to be executed by one thread as soon as all
 6552 threads reach the first barrier. The second barrier prevents the other threads from
 6553 proceeding until the serial section being executed by the one thread has completed.

6554 Although barriers can be implemented with mutexes and condition variables, the
 6555 referenced Almasi and Gottlieb (1989) provides ample illustration that such
 6556 implementations are significantly less efficient than is possible. While the relative
 6557 efficiency of barriers may well vary by implementation, it is important that they be
 6558 recognized in the IEEE Std 1003.1-200x to facilitate applications portability while providing
 6559 the necessary freedom to implementors.

- 6560 • Lack of Timeout Feature

6561 Alternate versions of most blocking routines have been provided to support watchdog
 6562 timeouts. No alternate interface of this sort has been provided for barrier waits for the
 6563 following reasons:

- 6564 • Multiple threads may use different timeout values, some of which may be indefinite.
 6565 It is not clear which threads should break through the barrier with a timeout error if
 6566 and when these timeouts expire.

- The barrier may become unusable once a thread breaks out of a `pthread_barrier_wait()` with a timeout error. There is, in general, no way to guarantee the consistency of a barrier's internal data structures once a thread has timed out of a `pthread_barrier_wait()`. Even the inclusion of a special barrier reinitialization function would not help much since it is not clear how this function would affect the behavior of threads that reach the barrier between the original timeout and the call to the reinitialization function.

Spin Locks

- Background

Spin locks represent an extremely low-level synchronization mechanism suitable primarily for use on shared memory multi-processors. It is typically an atomically modified Boolean value that is set to one when the lock is held and to zero when the lock is freed.

When a caller requests a spin lock that is already held, it typically spins in a loop testing whether the lock has become available. Such spinning wastes processor cycles so the lock should only be held for short durations and not across sleep/block operations. Callers should unlock spin locks before calling sleep operations.

Spin locks are available on a variety of systems. The functions included in IEEE Std 1003.1-200x are an attempt to standardize that existing practice.

- Lack of Timeout Feature

Alternate versions of most blocking routines have been provided to support watchdog timeouts. No alternate interface of this sort has been provided for spin locks for the following reasons:

- It is impossible to determine appropriate timeout intervals for spin locks in a portable manner. The amount of time one can expect to spend spin-waiting is inversely proportional to the degree of parallelism provided by the system.

It can vary from a few cycles when each competing thread is running on its own processor, to an indefinite amount of time when all threads are multiplexed on a single processor (which is why spin locking is not advisable on uniprocessors).

- When used properly, the amount of time the calling thread spends waiting on a spin lock should be considerably less than the time required to set up a corresponding watchdog timer. Since the primary purpose of spin locks is to provide a low-overhead synchronization mechanism for multi-processors, the overhead of a timeout mechanism was deemed unacceptable.

It was also suggested that an additional `count` argument be provided (on the `pthread_spin_lock()` call) in lieu of a true timeout so that a spin lock call could fail gracefully if it was unable to apply the lock after `count` attempts. This idea was rejected because it is not existing practice. Furthermore, the same effect can be obtained with `pthread_spin_trylock()`, as illustrated below:

```
int n = MAX_SPIN;

while ( --n >= 0 )
{
    if ( !pthread_spin_try_lock(...) )
        break;
}
if ( n >= 0 )
{
```

```

6613         /* Successfully acquired the lock */
6614     }
6615     else
6616     {
6617         /* Unable to acquire the lock */
6618     }

```

6619 • *process-shared* Attribute

6620 The initialization functions associated with most POSIX synchronization objects (for
6621 example, mutexes, barriers, and read-write locks) take an attributes object with a *process-*
6622 *shared* attribute that specifies whether or not the object is to be shared across processes. In
6623 the draft corresponding to the first balloting round, two separate initialization functions
6624 are provided for spin locks, however: one for spin locks that were to be shared across
6625 processes (*spin_init()*), and one for locks that were only used by multiple threads within a
6626 single process (*pthread_spin_init()*). This was done so as to keep the overhead associated
6627 with spin waiting to an absolute minimum. However, the balloting group requested that,
6628 since the overhead associated to a bit check was small, spin locks should be consistent with
6629 the rest of the synchronization primitives, and thus the *process-shared* attribute was
6630 introduced for spin locks.

6631 • Spin Locks *versus* Mutexes

6632 It has been suggested that mutexes are an adequate synchronization mechanism and spin
6633 locks are not necessary. Locking mechanisms typically must trade off the processor
6634 resources consumed while setting up to block the thread and the processor resources
6635 consumed by the thread while it is blocked. Spin locks require very little resources to set
6636 up the blocking of a thread. Existing practice is to simply loop, repeating the atomic
6637 locking operation until the lock is available. While the resources consumed to set up
6638 blocking of the thread are low, the thread continues to consume processor resources while
6639 it is waiting.

6640 On the other hand, mutexes may be implemented such that the processor resources
6641 consumed to block the thread are large relative to a spin lock. After detecting that the
6642 mutex lock is not available, the thread must alter its scheduling state, add itself to a set of
6643 waiting threads, and, when the lock becomes available again, undo all of this before taking
6644 over ownership of the mutex. However, while a thread is blocked by a mutex, no processor
6645 resources are consumed.

6646 Therefore, spin locks and mutexes may be implemented to have different characteristics.
6647 Spin locks may have lower overall overhead for very short-term blocking, and mutexes
6648 may have lower overall overhead when a thread will be blocked for longer periods of time.
6649 The presence of both interfaces allows implementations with these two different
6650 characteristics, both of which may be useful to a particular application.

6651 It has also been suggested that applications can build their own spin locks from the
6652 *pthread_mutex_trylock()* function:

```

6653 while (pthread_mutex_trylock(&mutex));

```

6654 The apparent simplicity of this construct is somewhat deceiving, however. While the actual
6655 wait is quite efficient, various guarantees on the integrity of mutex objects (for example,
6656 priority inheritance rules) may add overhead to the successful path of the trylock
6657 operation that is not required of spin locks. One could, of course, add an attribute to the
6658 mutex to bypass such overhead, but the very act of finding and testing this attribute
6659 represents more overhead than is found in the typical spin lock.

6660 The need to hold spin lock overhead to an absolute minimum also makes it impossible to

6661 provide guarantees against starvation similar to those provided for mutexes or read-write
6662 locks. The overhead required to implement such guarantees (for example, disabling
6663 preemption before spinning) may well exceed the overhead of the spin wait itself by many
6664 orders of magnitude. If a “safe” spin wait seems desirable, it can always be provided
6665 (albeit at some performance cost) via appropriate mutex attributes.

6666 **Robust Mutexes**

6667 Robust mutexes are intended to protect applications that use mutexes to protect data shared
6668 between different processes. If a process is terminated by a signal while a thread is holding a
6669 mutex, there is no chance for the process to clean up after it. Waiters for the locked mutex might
6670 wait indefinitely.

6671 With robust mutexes the problem can be solved: whenever a fatal signal terminates a process,
6672 current or future waiters of the mutex are notified about this fact. The locking function provides
6673 notification of this condition through the error condition [EOWNERDEAD]. A thread then has
6674 the chance to clean up the state protected by the mutex and mark the state as consistent again by
6675 a call to `pthread_mutex_consistent()`.

6676 Pre-existing implementations have used the semantics of robust mutexes for a variety of
6677 situations, some of them not defined in the standard. Where a normally terminated process (i.e.,
6678 when one thread calls `exit()`) causes notification of other waiters of robust mutexes if the mutex
6679 is locked by any thread in the process. This behavior is defined in the standard and makes sense
6680 because no thread other than the thread calling `exit()` has the chance to clean up its data.

6681 If a thread is terminated by cancellation or if it calls `pthread_exit()`, the situation is different. In
6682 both these situations the thread has the chance to clean up after itself by registering appropriate
6683 cleanup handlers. There is no real reason to demand that other waiters for a robust mutex the
6684 terminating thread owns are notified. The committee felt that this is actively encouraging bad
6685 practice because programmers are tempted to rely on the robust mutex semantics instead of
6686 correctly cleaning up after themselves.

6687 Therefore, the standard does not require notification of other waiters at the time a thread is
6688 terminated while the process continues to run. The mutex is still recognized as being locked by
6689 the process (with the thread gone it makes no sense to refer to the thread owning the mutex).
6690 Therefore, a terminating process will cause notifications about the dead owner to be sent to all
6691 waiters. This delay in the notification is not required, but programmers cannot rely on prompt
6692 notification after a thread is terminated.

6693 For the same reason is it not required that an implementation supports robust mutexes that are
6694 not shared between processes. If a robust mutex is used only within one process, all the cleanup
6695 can be performed by the threads themselves by registering appropriate cleanup handlers. Fatal
6696 signals are of no importance in this case because after the signal is delivered there is no thread
6697 remaining to use the mutex.

6698 Some implementations might choose to support intra-process robust mutexes and they might
6699 also send notification of a dead owner right after the previous owner died. But applications
6700 must not rely on this. Applications should only use robust mutexes for the purpose of handling
6701 fatal signals in situations where inter-process mutexes are in use.

6702 **Supported Threads Functions**

6703 On POSIX-conforming systems, the following symbolic constants are always conforming:

6704 `_POSIX_READER_WRITER_LOCKS`
6705 `_POSIX_THREADS`

6706 Therefore, the following threads functions are always supported:

6707	<code>pthread_atfork()</code>	<code>pthread_mutex_destroy()</code>
6708	<code>pthread_attr_destroy()</code>	<code>pthread_mutex_init()</code>
6709	<code>pthread_attr_getdetachstate()</code>	<code>pthread_mutex_lock()</code>
6710	<code>pthread_attr_getguardsize()</code>	<code>pthread_mutex_trylock()</code>
6711	<code>pthread_attr_getschedparam()</code>	<code>pthread_mutex_unlock()</code>
6712	<code>pthread_attr_init()</code>	<code>pthread_mutexattr_destroy()</code>
6713	<code>pthread_attr_setdetachstate()</code>	<code>pthread_mutexattr_getpshared()</code>
6714	<code>pthread_attr_setguardsize()</code>	<code>pthread_mutexattr_gettype()</code>
6715	<code>pthread_attr_setschedparam()</code>	<code>pthread_mutexattr_init()</code>
6716	<code>pthread_cancel()</code>	<code>pthread_mutexattr_setpshared()</code>
6717	<code>pthread_cleanup_pop()</code>	<code>pthread_mutexattr_settype()</code>
6718	<code>pthread_cleanup_push()</code>	<code>pthread_once()</code>
6719	<code>pthread_cond_broadcast()</code>	<code>pthread_rwlock_destroy()</code>
6720	<code>pthread_cond_destroy()</code>	<code>pthread_rwlock_init()</code>
6721	<code>pthread_cond_init()</code>	<code>pthread_rwlock_rdlock()</code>
6722	<code>pthread_cond_signal()</code>	<code>pthread_rwlock_tryrdlock()</code>
6723	<code>pthread_cond_timedwait()</code>	<code>pthread_rwlock_trywrlock()</code>
6724	<code>pthread_cond_wait()</code>	<code>pthread_rwlock_unlock()</code>
6725	<code>pthread_condattr_destroy()</code>	<code>pthread_rwlock_wrlock()</code>
6726	<code>pthread_condattr_getpshared()</code>	<code>pthread_rwlockattr_destroy()</code>
6727	<code>pthread_condattr_init()</code>	<code>pthread_rwlockattr_getpshared()</code>
6728	<code>pthread_condattr_setpshared()</code>	<code>pthread_rwlockattr_init()</code>
6729	<code>pthread_create()</code>	<code>pthread_rwlockattr_setpshared()</code>
6730	<code>pthread_detach()</code>	<code>pthread_self()</code>
6731	<code>pthread_equal()</code>	<code>pthread_setcancelstate()</code>
6732	<code>pthread_exit()</code>	<code>pthread_setcanceltype()</code>
6733	<code>pthread_getconcurrency()</code>	<code>pthread_setconcurrency()</code>
6734	<code>pthread_getspecific()</code>	<code>pthread_setspecific()</code>
6735	<code>pthread_join()</code>	<code>pthread_sigmask()</code>
6736	<code>pthread_key_create()</code>	<code>pthread_testcancel()</code>
6737	<code>pthread_key_delete()</code>	<code>sigwait()</code>
6738	<code>pthread_kill()</code>	

6739 On POSIX-conforming systems, the symbolic constant `_POSIX_THREAD_SAFE_FUNCTIONS` is
6740 always defined. Therefore, the following functions are always supported:

6741	<code>asctime_r()</code>	<code>getpwuid_r()</code>
6742	<code>ctime_r()</code>	<code>gmtime_r()</code>
6743	<code>flockfile()</code>	<code>localtime_r()</code>
6744	<code>ftrylockfile()</code>	<code>putc_unlocked()</code>
6745	<code>funlockfile()</code>	<code>putchar_unlocked()</code>
6746	<code>getc_unlocked()</code>	<code>rand_r()</code>
6747	<code>getchar_unlocked()</code>	<code>readdir_r()</code>
6748	<code>getgrgid_r()</code>	<code>strerror_r()</code>
6749	<code>getgrnam_r()</code>	<code>strtok_r()</code>
6750	<code>getpwnam_r()</code>	

6751 Threads Extensions

6752 The following extensions to the IEEE P1003.1c draft standard are now supported in
6753 IEEE Std 1003.1-200x as part of the alignment with the Single UNIX Specification:

- 6754 • Extended mutex attribute types
- 6755 • Read-write locks and attributes (also introduced by the IEEE Std 1003.1j-2000 amendment)
- 6756 • Thread concurrency level
- 6757 • Thread stack guard size
- 6758 • Parallel I/O
- 6759 • Robust mutexes

6760 These extensions carefully follow the threads programming model specified in POSIX.1c. As
6761 with POSIX.1c, all the new functions return zero if successful; otherwise, an error number is
6762 returned to indicate the error.

6763 The concept of attribute objects was introduced in POSIX.1c to allow implementations to extend
6764 IEEE Std 1003.1-200x without changing the existing interfaces. Attribute objects were defined for
6765 threads, mutexes, and condition variables. Attributes objects are defined as implementation-
6766 defined opaque types to aid extensibility, and functions are defined to allow attributes to be set
6767 or retrieved. This model has been followed when adding the new type attribute of
6768 `pthread_mutexattr_t` or the new read-write lock attributes object `pthread_rwlockattr_t`.

- 6769 • Extended Mutex Attributes

6770 POSIX.1c defines a mutex attributes object as an implementation-defined opaque object of
6771 type `pthread_mutexattr_t`, and specifies a number of attributes which this object must
6772 have and a number of functions which manipulate these attributes. These attributes
6773 include *detachstate*, *inheritsched*, *schedparm*, *schedpolicy*, *contentionscope*, *stackaddr*, and
6774 *stacksize*.

6775 The System Interfaces volume of IEEE Std 1003.1-200x specifies another mutex attribute
6776 called *type*. The *type* attribute allows applications to specify the behavior of mutex locking
6777 operations in situations where POSIX.1c behavior is undefined. The OSF DCE threads
6778 implementation, based on Draft 4 of POSIX.1c, specified a similar attribute. Note that the
6779 names of the attributes have changed somewhat from the OSF DCE threads
6780 implementation.

6781 The System Interfaces volume of IEEE Std 1003.1-200x also extends the specification of the
6782 following POSIX.1c functions which manipulate mutexes:

6783 `pthread_mutex_lock()`
 6784 `pthread_mutex_trylock()`
 6785 `pthread_mutex_unlock()`

6786 to take account of the new mutex attribute type and to specify behavior which was
 6787 declared as undefined in POSIX.1c. How a calling thread acquires or releases a mutex now
 6788 depends upon the mutex *type* attribute.

6789 The *type* attribute can have the following values:

6790 PTHREAD_MUTEX_NORMAL

6791 Basic mutex with no specific error checking built in. Does not report a deadlock error.

6792 PTHREAD_MUTEX_RECURSIVE

6793 Allows any thread to recursively lock a mutex. The mutex must be unlocked an equal
 6794 number of times to release the mutex.

6795 PTHREAD_MUTEX_ERRORCHECK

6796 Detects and reports simple usage errors; that is, an attempt to unlock a mutex that is
 6797 not locked by the calling thread or that is not locked at all, or an attempt to relock a
 6798 mutex the thread already owns.

6799 PTHREAD_MUTEX_DEFAULT

6800 The default mutex type. May be mapped to any of the above mutex types or may be
 6801 an implementation-defined type.

6802 *Normal* mutexes do not detect deadlock conditions; for example, a thread will hang if it
 6803 tries to relock a normal mutex that it already owns. Attempting to unlock a mutex locked
 6804 by another thread, or unlocking an unlocked mutex, results in undefined behavior. Normal
 6805 mutexes will usually be the fastest type of mutex available on a platform but provide the
 6806 least error checking.

6807 *Recursive* mutexes are useful for converting old code where it is difficult to establish clear
 6808 boundaries of synchronization. A thread can relock a recursive mutex without first
 6809 unlocking it. The relocking deadlock which can occur with normal mutexes cannot occur
 6810 with this type of mutex. However, multiple locks of a recursive mutex require the same
 6811 number of unlocks to release the mutex before another thread can acquire the mutex.
 6812 Furthermore, this type of mutex maintains the concept of an owner. Thus, a thread
 6813 attempting to unlock a recursive mutex which another thread has locked returns with an
 6814 error. A thread attempting to unlock a recursive mutex that is not locked returns with an
 6815 error. Never use a recursive mutex with condition variables because the implicit unlock
 6816 performed by `pthread_cond_wait()` or `pthread_cond_timedwait()` will not actually release the
 6817 mutex if it had been locked multiple times.

6818 *Errorcheck* mutexes provide error checking and are useful primarily as a debugging aid. A
 6819 thread attempting to relock an errorcheck mutex without first unlocking it returns with an
 6820 error. Again, this type of mutex maintains the concept of an owner. Thus, a thread
 6821 attempting to unlock an errorcheck mutex which another thread has locked returns with
 6822 an error. A thread attempting to unlock an errorcheck mutex that is not locked also returns
 6823 with an error. It should be noted that errorcheck mutexes will almost always be much
 6824 slower than normal mutexes due to the extra state checks performed.

6825 The default mutex type provides implementation-defined error checking. The default
 6826 mutex may be mapped to one of the other defined types or may be something entirely
 6827 different. This enables each vendor to provide the mutex semantics which the vendor feels
 6828 will be most useful to their target users. Most vendors will probably choose to make
 6829 normal mutexes the default so as to give applications the benefit of the fastest type of

6830 mutexes available on their platform. Check your implementation's documentation.

6831 An application developer can use any of the mutex types almost interchangeably as long
 6832 as the application does not depend upon the implementation detecting (or failing to
 6833 detect) any particular errors. Note that a recursive mutex can be used with condition
 6834 variable waits as long as the application never recursively locks the mutex.

6835 Two functions are provided for manipulating the *type* attribute of a mutex attributes object.
 6836 This attribute is set or returned in the *type* parameter of these functions. The
 6837 `pthread_mutexattr_settype()` function is used to set a specific type value while
 6838 `pthread_mutexattr_gettype()` is used to return the type of the mutex. Setting the *type*
 6839 attribute of a mutex attributes object affects only mutexes initialized using that mutex
 6840 attributes object. Changing the *type* attribute does not affect mutexes previously initialized
 6841 using that mutex attributes object.

6842 • Read-Write Locks and Attributes

6843 The read-write locks introduced have been harmonized with those in IEEE Std
 6844 1003.1j-2000; see also [Section B.2.9.6](#) (on page 175).

6845 Read-write locks (also known as reader-writer locks) allow a thread to exclusively lock
 6846 some shared data while updating that data, or allow any number of threads to have
 6847 simultaneous read-only access to the data.

6848 Unlike a mutex, a read-write lock distinguishes between reading data and writing data. A
 6849 mutex excludes all other threads. A read-write lock allows other threads access to the data,
 6850 providing no thread is modifying the data. Thus, a read-write lock is less primitive than
 6851 either a mutex-condition variable pair or a semaphore.

6852 Application developers should consider using a read-write lock rather than a mutex to
 6853 protect data that is frequently referenced but seldom modified. Most threads (readers) will
 6854 be able to read the data without waiting and will only have to block when some other
 6855 thread (a writer) is in the process of modifying the data. Conversely a thread that wants to
 6856 change the data is forced to wait until there are no readers. This type of lock is often used
 6857 to facilitate parallel access to data on multi-processor platforms or to avoid context
 6858 switches on single processor platforms where multiple threads access the same data.

6859 If a read-write lock becomes unlocked and there are multiple threads waiting to acquire
 6860 the write lock, the implementation's scheduling policy determines which thread acquires
 6861 the read-write lock for writing. If there are multiple threads blocked on a read-write lock
 6862 for both read locks and write locks, it is unspecified whether the readers or a writer
 6863 acquire the lock first. However, for performance reasons, implementations often favor
 6864 writers over readers to avoid potential writer starvation.

6865 A read-write lock object is an implementation-defined opaque object of type
 6866 `pthread_rwlock_t` as defined in `<pthread.h>`. There are two different sorts of locks
 6867 associated with a read-write lock: a read lock and a write lock.

6868 The `pthread_rwlockattr_init()` function initializes a read-write lock attributes object with the
 6869 default value for all the attributes defined in the implementation. After a read-write lock
 6870 attributes object has been used to initialize one or more read-write locks, changes to the
 6871 read-write lock attributes object, including destruction, do not affect previously initialized
 6872 read-write locks.

6873 Implementations must provide at least the read-write lock attribute *process-shared*. This
 6874 attribute can have the following values:

PTHREAD_PROCESS_SHARED

Any thread of any process that has access to the memory where the read-write lock resides can manipulate the read-write lock.

PTHREAD_PROCESS_PRIVATE

Only threads created within the same process as the thread that initialized the read-write lock can manipulate the read-write lock. This is the default value.

The `pthread_rwlockattr_setpshared()` function is used to set the *process-shared* attribute of an initialized read-write lock attributes object while the function `pthread_rwlockattr_getpshared()` obtains the current value of the *process-shared* attribute.

A read-write lock attributes object is destroyed using the `pthread_rwlockattr_destroy()` function. The effect of subsequent use of the read-write lock attributes object is undefined.

A thread creates a read-write lock using the `pthread_rwlock_init()` function. The attributes of the read-write lock can be specified by the application developer; otherwise, the default implementation-defined read-write lock attributes are used if the pointer to the read-write lock attributes object is NULL. In cases where the default attributes are appropriate, the PTHREAD_RWLOCK_INITIALIZER macro can be used to initialize statically allocated read-write locks.

A thread which wants to apply a read lock to the read-write lock can use either `pthread_rwlock_rdlock()` or `pthread_rwlock_tryrdlock()`. If `pthread_rwlock_rdlock()` is used, the thread acquires a read lock if a writer does not hold the write lock and there are no writers blocked on the write lock. If a read lock is not acquired, the calling thread blocks until it can acquire a lock. However, if `pthread_rwlock_tryrdlock()` is used, the function returns immediately with the error [EBUSY] if any thread holds a write lock or there are blocked writers waiting for the write lock.

A thread which wants to apply a write lock to the read-write lock can use either of two functions: `pthread_rwlock_wrlock()` or `pthread_rwlock_trywrlock()`. If `pthread_rwlock_wrlock()` is used, the thread acquires the write lock if no other reader or writer threads hold the read-write lock. If the write lock is not acquired, the thread blocks until it can acquire the write lock. However, if `pthread_rwlock_trywrlock()` is used, the function returns immediately with the error [EBUSY] if any thread is holding either a read or a write lock.

The `pthread_rwlock_unlock()` function is used to unlock a read-write lock object held by the calling thread. Results are undefined if the read-write lock is not held by the calling thread. If there are other read locks currently held on the read-write lock object, the read-write lock object remains in the read locked state but without the current thread as one of its owners. If this function releases the last read lock for this read-write lock object, the read-write lock object is put in the unlocked read state. If this function is called to release a write lock for this read-write lock object, the read-write lock object is put in the unlocked state.

- Thread Concurrency Level

On threads implementations that multiplex user threads onto a smaller set of kernel execution entities, the system attempts to create a reasonable number of kernel execution entities for the application upon application startup.

On some implementations, these kernel entities are retained by user threads that block in the kernel. Other implementations do not *timeslice* user threads so that multiple compute-bound user threads can share a kernel thread. On such implementations, some applications may use up all the available kernel execution entities before their user-space threads are used up. The process may be left with user threads capable of doing work for the application but with no way to schedule them.

6922 The *pthread_setconcurrency()* function enables an application to request more kernel
 6923 entities; that is, specify a desired concurrency level. However, this function merely
 6924 provides a hint to the implementation. The implementation is free to ignore this request or
 6925 to provide some other number of kernel entities. If an implementation does not multiplex
 6926 user threads onto a smaller number of kernel execution entities, the
 6927 *pthread_setconcurrency()* function has no effect.

6928 The *pthread_setconcurrency()* function may also have an effect on implementations where
 6929 the kernel mode and user mode schedulers cooperate to ensure that ready user threads are
 6930 not prevented from running by other threads blocked in the kernel.

6931 The *pthread_getconcurrency()* function always returns the value set by a previous call to
 6932 *pthread_setconcurrency()*. However, if *pthread_setconcurrency()* was not previously called,
 6933 this function returns zero to indicate that the threads implementation is maintaining the
 6934 concurrency level.

6935 • Thread Stack Guard Size

6936 DCE threads introduced the concept of a “thread stack guard size”. Most thread
 6937 implementations add a region of protected memory to a thread’s stack, commonly known
 6938 as a “guard region”, as a safety measure to prevent stack pointer overflow in one thread
 6939 from corrupting the contents of another thread’s stack. The default size of the guard
 6940 regions attribute is {PAGESIZE} bytes and is implementation-defined.

6941 Some application developers may wish to change the stack guard size. When an
 6942 application creates a large number of threads, the extra page allocated for each stack may
 6943 strain system resources. In addition to the extra page of memory, the kernel’s memory
 6944 manager has to keep track of the different protections on adjoining pages. When this is a
 6945 problem, the application developer may request a guard size of 0 bytes to conserve system
 6946 resources by eliminating stack overflow protection.

6947 Conversely an application that allocates large data structures such as arrays on the stack
 6948 may wish to increase the default guard size in order to detect stack overflow. If a thread
 6949 allocates two pages for a data array, a single guard page provides little protection against
 6950 thread stack overflows since the thread can corrupt adjoining memory beyond the guard
 6951 page.

6952 The System Interfaces volume of IEEE Std 1003.1-200x defines a new attribute of a thread
 6953 attributes object; that is, the *guardsize* attribute which allows applications to specify the size
 6954 of the guard region of a thread’s stack.

6955 Two functions are provided for manipulating a thread’s stack guard size. The
 6956 *pthread_attr_setguardsize()* function sets the thread *guardsize* attribute, and the
 6957 *pthread_attr_getguardsize()* function retrieves the current value.

6958 An implementation may round up the requested guard size to a multiple of the
 6959 configurable system variable {PAGESIZE}. In this case, *pthread_attr_getguardsize()* returns
 6960 the guard size specified by the previous *pthread_attr_setguardsize()* function call and not
 6961 the rounded up value.

6962 If an application is managing its own thread stacks using the *stackaddr* attribute, the
 6963 *guardsize* attribute is ignored and no stack overflow protection is provided. In this case, it is
 6964 the responsibility of the application to manage stack overflow along with stack allocation.

6965 • Parallel I/O

6966 Suppose two or more threads independently issue read requests on the same file. To read
 6967 specific data from a file, a thread must first call *lseek()* to seek to the proper offset in the
 6968 file, and then call *read()* to retrieve the required data. If more than one thread does this at

6969 the same time, the first thread may complete its seek call, but before it gets a chance to
 6970 issue its read call a second thread may complete its seek call, resulting in the first thread
 6971 accessing incorrect data when it issues its read call. One workaround is to lock the file
 6972 descriptor while seeking and reading or writing, but this reduces parallelism and adds
 6973 overhead.

6974 Instead, the System Interfaces volume of IEEE Std 1003.1-200x provides two functions to
 6975 make seek/read and seek/write operations atomic. The file descriptor's current offset is
 6976 unchanged, thus allowing multiple read and write operations to proceed in parallel. This
 6977 improves the I/O performance of threaded applications. The *pread()* function is used to do
 6978 an atomic read of data from a file into a buffer. Conversely, the *pwrite()* function does an
 6979 atomic write of data from a buffer to a file.

6980 B.2.9.1 Thread-Safety

6981 All functions required by IEEE Std 1003.1-200x need to be thread-safe. Implementations have to
 6982 provide internal synchronization when necessary in order to achieve this goal. In certain cases—
 6983 for example, most floating-point implementations—context switch code may have to manage
 6984 the writable shared state.

6985 While a read from a pipe of $\{\text{PIPE_MAX}\} \times 2$ bytes may not generate a single atomic and thread-
 6986 safe stream of bytes, it should generate “several” (individually atomic) thread-safe streams of
 6987 bytes. Similarly, while reading from a terminal device may not generate a single atomic and
 6988 thread-safe stream of bytes, it should generate some finite number of (individually atomic) and
 6989 thread-safe streams of bytes. That is, concurrent calls to read for a pipe, FIFO, or terminal device
 6990 are not allowed to result in corrupting the stream of bytes or other internal data. However,
 6991 *read()*, in these cases, is not required to return a single contiguous and atomic stream of bytes.

6992 It is not required that all functions provided by IEEE Std 1003.1-200x be either async-cancel-safe
 6993 or async-signal-safe.

6994 As it turns out, some functions are inherently not thread-safe; that is, their interface
 6995 specifications preclude reentrancy. For example, some functions (such as *asctime()*) return a
 6996 pointer to a result stored in memory space allocated by the function on a per-process basis. Such
 6997 a function is not thread-safe, because its result can be overwritten by successive invocations.
 6998 Other functions, while not inherently non-thread-safe, may be implemented in ways that lead to
 6999 them not being thread-safe. For example, some functions (such as *rand()*) store state information
 7000 (such as a seed value, which survives multiple function invocations) in memory space allocated
 7001 by the function on a per-process basis. The implementation of such a function is not thread-safe
 7002 if the implementation fails to synchronize invocations of the function and thus fails to protect
 7003 the state information. The problem is that when the state information is not protected,
 7004 concurrent invocations can interfere with one another (for example, applications using *rand()*
 7005 may see the same seed value).

7006 Thread-Safety and Locking of Existing Functions

7007 Originally, POSIX.1 was not designed to work in a multi-threaded environment, and some
 7008 implementations of some existing functions will not work properly when executed concurrently.
 7009 To provide routines that will work correctly in an environment with threads (“thread-safe”), two
 7010 problems need to be solved:

- 7011 1. Routines that maintain or return pointers to static areas internal to the routine (which
 7012 may now be shared) need to be modified. The routines *ttyname()* and *localtime()* are
 7013 examples.
- 7014 2. Routines that access data space shared by more than one thread need to be modified. The
 7015 *malloc()* function and the *stdio* family routines are examples.

7016 There are a variety of constraints on these changes. The first is compatibility with the existing
 7017 versions of these functions—non-thread-safe functions will continue to be in use for some time,
 7018 as the original interfaces are used by existing code. Another is that the new thread-safe versions
 7019 of these functions represent as small a change as possible over the familiar interfaces provided
 7020 by the existing non-thread-safe versions. The new interfaces should be independent of any
 7021 particular threads implementation. In particular, they should be thread-safe without depending
 7022 on explicit thread-specific memory. Finally, there should be minimal performance penalty due to
 7023 the changes made to the functions.

7024 It is intended that the list of functions from POSIX.1 that cannot be made thread-safe and for
 7025 which corrected versions are provided be complete.

7026 *Thread-Safety and Locking Solutions*

7027 Many of the POSIX.1 functions were thread-safe and did not change at all. However, some
 7028 functions (for example, the math functions typically found in **libm**) are not thread-safe because
 7029 of writable shared global state. For instance, in IEEE Std 754-1985 floating-point
 7030 implementations, the computation modes and flags are global and shared.

7031 Some functions are not thread-safe because a particular implementation is not reentrant,
 7032 typically because of a non-essential use of static storage. These require only a new
 7033 implementation.

7034 Thread-safe libraries are useful in a wide range of parallel (and asynchronous) programming
 7035 environments, not just within pthreads. In order to be used outside the context of pthreads,
 7036 however, such libraries still have to use some synchronization method. These could either be
 7037 independent of the pthread synchronization operations, or they could be a subset of the pthread
 7038 interfaces. Either method results in thread-safe library implementations that can be used without
 7039 the rest of pthreads.

7040 Some functions, such as the *stdio* family interface and dynamic memory allocation functions
 7041 such as *malloc()*, are inter-dependent routines that share resources (for example, buffers) across
 7042 related calls. These require synchronization to work correctly, but they do not require any
 7043 change to their external (user-visible) interfaces.

7044 In some cases, such as *getc()* and *putc()*, adding synchronization is likely to create an
 7045 unacceptable performance impact. In this case, slower thread-safe synchronized functions are to
 7046 be provided, but the original, faster (but unsafe) functions (which may be implemented as
 7047 macros) are retained under new names. Some additional special-purpose synchronization
 7048 facilities are necessary for these macros to be usable in multi-threaded programs. This also
 7049 requires changes in **<stdio.h>**.

7050 The other common reason that functions are unsafe is that they return a pointer to static storage,
 7051 making the functions non-thread-safe. This has to be changed, and there are three natural
 7052 choices:

7053 1. Return a pointer to thread-specific storage

7054 This could incur a severe performance penalty on those architectures with a costly
 7055 implementation of the thread-specific data interface.

7056 A variation on this technique is to use *malloc()* to allocate storage for the function output
 7057 and return a pointer to this storage. This technique may also have an undesirable
 7058 performance impact, however, and a simplistic implementation requires that the user
 7059 program explicitly free the storage object when it is no longer needed. This technique is
 7060 used by some existing POSIX.1 functions. With careful implementation for infrequently
 7061 used functions, there may be little or no performance or storage penalty, and the
 7062 maintenance of already-standardized interfaces is a significant benefit.

7063 2. Return the actual value computed by the function

7064 This technique can only be used with functions that return pointers to structures—
 7065 routines that return character strings would have to wrap their output in an enclosing
 7066 structure in order to return the output on the stack. There is also a negative performance
 7067 impact inherent in this solution in that the output value has to be copied twice before it
 7068 can be used by the calling function: once from the called routine's local buffers to the top
 7069 of the stack, then from the top of the stack to the assignment target. Finally, many older
 7070 compilers cannot support this technique due to a historical tendency to use internal static
 7071 buffers to deliver the results of structure-valued functions.

7072 3. Have the caller pass the address of a buffer to contain the computed value

7073 The only disadvantage of this approach is that extra arguments have to be provided by
 7074 the calling program. It represents the most efficient solution to the problem, however,
 7075 and, unlike the *malloc()* technique, it is semantically clear.

7076 There are some routines (often groups of related routines) whose interfaces are inherently non-
 7077 thread-safe because they communicate across multiple function invocations by means of static
 7078 memory locations. The solution is to redesign the calls so that they are thread-safe, typically by
 7079 passing the needed data as extra parameters. Unfortunately, this may require major changes to
 7080 the interface as well.

7081 A floating-point implementation using IEEE Std 754-1985 is a case in point. A less problematic
 7082 example is the *rand48* family of pseudo-random number generators. The functions *getgrgid()*,
 7083 *getgrnam()*, *getpwnam()*, and *getpwuid()* are another such case.

7084 The problems with *errno* are discussed in [Alternative Solutions for Per-Thread *errno*](#) (on page
 7085 91).

7086 Some functions can be thread-safe or not, depending on their arguments. These include the
 7087 *tmpnam()* and *ctermid()* functions. These functions have pointers to character strings as
 7088 arguments. If the pointers are not NULL, the functions store their results in the character string;
 7089 however, if the pointers are NULL, the functions store their results in an area that may be static
 7090 and thus subject to overwriting by successive calls. These should only be called by multi-thread
 7091 applications when their arguments are non-NULL.

7092 *Asynchronous Safety and Thread-Safety*

7093 A floating-point implementation has many modes that effect rounding and other aspects of
 7094 computation. Functions in some math library implementations may change the computation
 7095 modes for the duration of a function call. If such a function call is interrupted by a signal or
 7096 cancellation, the floating-point state is not required to be protected.

7097 There is a significant cost to make floating-point operations async-cancel-safe or async-signal-
 7098 safe; accordingly, neither form of async safety is required.

7099 *Functions Returning Pointers to Static Storage*

7100 For those functions that are not thread-safe because they return values in fixed size statically
 7101 allocated structures, alternate “_r” forms are provided that pass a pointer to an explicit result
 7102 structure. Those that return pointers into library-allocated buffers have forms provided with
 7103 explicit buffer and length parameters.

7104 For functions that return pointers to library-allocated buffers, it makes sense to provide “_r”
 7105 versions that allow the application control over allocation of the storage in which results are
 7106 returned. This allows the state used by these functions to be managed on an application-specific
 7107 basis, supporting per-thread, per-process, or other application-specific sharing relationships.

7108 Early proposals had provided “_r” versions for functions that returned pointers to variable-size

7109 buffers without providing a means for determining the required buffer size. This would have
 7110 made using such functions exceedingly clumsy, potentially requiring iteratively calling them
 7111 with increasingly larger guesses for the amount of storage required. Hence, *sysconf()* variables
 7112 have been provided for such functions that return the maximum required buffer size.

7113 Thus, the rule that has been followed by IEEE Std 1003.1-200x when adapting single-threaded
 7114 non-thread-safe functions is as follows: all functions returning pointers to library-allocated
 7115 storage should have “_r” versions provided, allowing the application control over the storage
 7116 allocation. Those with variable-sized return values accept both a buffer address and a length
 7117 parameter. The *sysconf()* variables are provided to supply the appropriate buffer sizes when
 7118 required. Implementors are encouraged to apply the same rule when adapting their own
 7119 existing functions to a pthreads environment.

7120 B.2.9.2 Thread IDs

7121 Separate applications should communicate through well-defined interfaces and should not
 7122 depend on each other’s implementation. For example, if a programmer decides to rewrite the
 7123 *sort* utility using multiple threads, it should be easy to do this so that the interface to the *sort*
 7124 utility does not change. Consider that if the user causes SIGINT to be generated while the *sort*
 7125 utility is running, keeping the same interface means that the entire *sort* utility is killed, not just
 7126 one of its threads. As another example, consider a realtime application that manages a reactor.
 7127 Such an application may wish to allow other applications to control the priority at which it
 7128 watches the control rods. One technique to accomplish this is to write the ID of the thread
 7129 watching the control rods into a file and allow other programs to change the priority of that
 7130 thread as they see fit. A simpler technique is to have the reactor process accept IPCs
 7131 (Interprocess Communication messages) from other processes, telling it at a semantic level what
 7132 priority the program should assign to watching the control rods. This allows the programmer
 7133 greater flexibility in the implementation. For example, the programmer can change the
 7134 implementation from having one thread per rod to having one thread watching all of the rods
 7135 without changing the interface. Having threads live inside the process means that the
 7136 implementation of a process is invisible to outside processes (excepting debuggers and system
 7137 management tools).

7138 Threads do not provide a protection boundary. Every thread model allows threads to share
 7139 memory with other threads and encourages this sharing to be widespread. This means that one
 7140 thread can wipe out memory that is needed for the correct functioning of other threads that are
 7141 sharing its memory. Consequently, providing each thread with its own user and/or group IDs
 7142 would not provide a protection boundary between threads sharing memory.

7143 B.2.9.3 Thread Mutexes

7144 There is no additional rationale provided for this section.

7145 B.2.9.4 Thread Scheduling

7146 • Scheduling Implementation Models

7147 The following scheduling implementation models are presented in terms of threads and
 7148 “kernel entities”. This is to simplify exposition of the models, and it does not imply that
 7149 an implementation actually has an identifiable “kernel entity”.

7150 A kernel entity is not defined beyond the fact that it has scheduling attributes that are used
 7151 to resolve contention with other kernel entities for execution resources. A kernel entity
 7152 may be thought of as an envelope that holds a thread or a separate kernel thread. It is not a
 7153 conventional process, although it shares with the process the attribute that it has a single
 7154 thread of control; it does not necessarily imply an address space, open files, and so on. It is

7155 better thought of as a primitive facility upon which conventional processes and threads
7156 may be constructed.

7157 — System Thread Scheduling Model

7158 This model consists of one thread per kernel entity. The kernel entity is solely
7159 responsible for scheduling thread execution on one or more processors. This model
7160 schedules all threads against all other threads in the system using the scheduling
7161 attributes of the thread.

7162 — Process Scheduling Model

7163 A generalized process scheduling model consists of two levels of scheduling. A
7164 threads library creates a pool of kernel entities, as required, and schedules threads to
7165 run on them using the scheduling attributes of the threads. Typically, the size of the
7166 pool is a function of the simultaneously runnable threads, not the total number of
7167 threads. The kernel then schedules the kernel entities onto processors according to
7168 their scheduling attributes, which are managed by the threads library. This set model
7169 potentially allows a wide range of mappings between threads and kernel entities.

7170 • System and Process Scheduling Model Performance

7171 There are a number of important implications on the performance of applications using
7172 these scheduling models. The process scheduling model potentially provides lower
7173 overhead for making scheduling decisions, since there is no need to access kernel-level
7174 information or functions and the set of schedulable entities is smaller (only the threads
7175 within the process).

7176 On the other hand, since the kernel is also making scheduling decisions regarding the
7177 system resources under its control (for example, CPU(s), I/O devices, memory), decisions
7178 that do not take thread scheduling parameters into account can result in unspecified
7179 delays for realtime application threads, causing them to miss maximum response time
7180 limits.

7181 • Rate Monotonic Scheduling

7182 Rate monotonic scheduling was considered, but rejected for standardization in the context
7183 of pthreads. A sporadic server policy is included.

7184 • Scheduling Options

7185 In IEEE Std 1003.1-200x, the basic thread scheduling functions are defined under the
7186 threads functionality, so that they are required of all threads implementations. However,
7187 there are no specific scheduling policies required by this functionality to allow for
7188 conforming thread implementations that are not targeted to realtime applications.

7189 Specific standard scheduling policies are defined to be under the Thread Execution
7190 Scheduling option, and they are specifically designed to support realtime applications by
7191 providing predictable resource-sharing sequences. The name of this option was chosen to
7192 emphasize that this functionality is defined as appropriate for realtime applications that
7193 require simple priority-based scheduling.

7194 It is recognized that these policies are not necessarily satisfactory for some multi-processor
7195 implementations, and work is ongoing to address a wider range of scheduling behaviors.
7196 The interfaces have been chosen to create abundant opportunity for future scheduling
7197 policies to be implemented and standardized based on this interface. In order to
7198 standardize a new scheduling policy, all that is required (from the standpoint of thread
7199 scheduling attributes) is to define a new policy name, new members of the thread
7200 attributes object, and functions to set these members when the scheduling policy is equal

7201 to the new value.

7202 **Scheduling Contention Scope**

7203 In order to accommodate the requirement for realtime response, each thread has a scheduling
7204 contention scope attribute. Threads with a system scheduling contention scope have to be
7205 scheduled with respect to all other threads in the system. These threads are usually bound to a
7206 single kernel entity that reflects their scheduling attributes and are directly scheduled by the
7207 kernel.

7208 Threads with a process scheduling contention scope need be scheduled only with respect to the
7209 other threads in the process. These threads may be scheduled within the process onto a pool of
7210 kernel entities. The implementation is also free to bind these threads directly to kernel entities
7211 and let them be scheduled by the kernel. Process scheduling contention scope allows the
7212 implementation the most flexibility and is the default if both contention scopes are supported
7213 and none is specified.

7214 Thus, the choice by implementors to provide one or the other (or both) of these scheduling
7215 models is driven by the need of their supported application domains for worst-case (that is,
7216 realtime) response, or average-case (non-realtime) response.

7217 **Scheduling Allocation Domain**

7218 The SCHED_FIFO and SCHED_RR scheduling policies take on different characteristics on a
7219 multi-processor. Other scheduling policies are also subject to changed behavior when executed
7220 on a multi-processor. The concept of scheduling allocation domain determines the set of
7221 processors on which the threads of an application may run. By considering the application's
7222 processor scheduling allocation domain for its threads, scheduling policies can be defined in
7223 terms of their behavior for varying processor scheduling allocation domain values. It is
7224 conceivable that not all scheduling allocation domain sizes make sense for all scheduling
7225 policies on all implementations. The concept of scheduling allocation domain, however, is a
7226 useful tool for the description of multi-processor scheduling policies.

7227 The "process control" approach to scheduling obtains significant performance advantages from
7228 dynamic scheduling allocation domain sizes when it is applicable.

7229 Non-Uniform Memory Access (NUMA) multi-processors may use a system scheduling structure
7230 that involves reassignment of threads among scheduling allocation domains. In NUMA
7231 machines, a natural model of scheduling is to match scheduling allocation domains to clusters of
7232 processors. Load balancing in such an environment requires changing the scheduling allocation
7233 domain to which a thread is assigned.

7234 **Scheduling Documentation**

7235 Implementation-provided scheduling policies need to be completely documented in order to be
7236 useful. This documentation includes a description of the attributes required for the policy, the
7237 scheduling interaction of threads running under this policy and all other supported policies, and
7238 the effects of all possible values for processor scheduling allocation domain. Note that for the
7239 implementor wishing to be minimally-compliant, it is (minimally) acceptable to define the
7240 behavior as undefined.

7241 Scheduling Contention Scope Attribute

7242 The scheduling contention scope defines how threads compete for resources. Within
7243 IEEE Std 1003.1-200x, scheduling contention scope is used to describe only how threads are
7244 scheduled in relation to one another in the system. That is, either they are scheduled against all
7245 other threads in the system (“system scope”) or only against those threads in the process
7246 (“process scope”). In fact, scheduling contention scope may apply to additional resources,
7247 including virtual timers and profiling, which are not currently considered by
7248 IEEE Std 1003.1-200x.

7249 Mixed Scopes

7250 If only one scheduling contention scope is supported, the scheduling decision is straightforward.
7251 To perform the processor scheduling decision in a mixed scope environment, it is necessary to
7252 map the scheduling attributes of the thread with process-wide contention scope to the same
7253 attribute space as the thread with system-wide contention scope.

7254 Since a conforming implementation has to support one and may support both scopes, it is useful
7255 to discuss the effects of such choices with respect to example applications. If an implementation
7256 supports both scopes, mixing scopes provides a means of better managing system-level (that is,
7257 kernel-level) and library-level resources. In general, threads with system scope will require the
7258 resources of a separate kernel entity in order to guarantee the scheduling semantics. On the
7259 other hand, threads with process scope can share the resources of a kernel entity while
7260 maintaining the scheduling semantics.

7261 The application is free to create threads with dedicated kernel resources, and other threads that
7262 multiplex kernel resources. Consider the example of a window server. The server allocates two
7263 threads per widget: one thread manages the widget user interface (including drawing), while
7264 the other thread takes any required application action. This allows the widget to be “active”
7265 while the application is computing. A screen image may be built from thousands of widgets. If
7266 each of these threads had been created with system scope, then most of the kernel-level
7267 resources might be wasted, since only a few widgets are active at any one time. In addition,
7268 mixed scope is particularly useful in a window server where one thread with high priority and
7269 system scope handles the mouse so that it tracks well. As another example, consider a database
7270 server. For each of the hundreds or thousands of clients supported by a large server, an
7271 equivalent number of threads will have to be created. If each of these threads were system scope,
7272 the consequences would be the same as for the window server example above. However, the
7273 server could be constructed so that actual retrieval of data is done by several dedicated threads.
7274 Dedicated threads that do work for all clients frequently justify the added expense of system
7275 scope. If it were not permissible to mix system and process threads in the same process, this type
7276 of solution would not be possible.

7277 Dynamic Thread Scheduling Parameters Access

7278 In many time-constrained applications, there is no need to change the scheduling attributes
7279 dynamically during thread or process execution, since the general use of these attributes is to
7280 reflect directly the time constraints of the application. Since these time constraints are generally
7281 imposed to meet higher-level system requirements, such as accuracy or availability, they
7282 frequently should remain unchanged during application execution.

7283 However, there are important situations in which the scheduling attributes should be changed.
7284 Generally, this will occur when external environmental conditions exist in which the time
7285 constraints change. Consider, for example, a space vehicle major mode change, such as the
7286 change from ascent to descent mode, or the change from the space environment to the
7287 atmospheric environment. In such cases, the frequency with which many of the sensors or
7288 actuators need to be read or written will change, which will necessitate a priority change. In

7289 other cases, even the existence of a time constraint might be temporary, necessitating not just a
 7290 priority change, but also a policy change for ongoing threads or processes. For this reason, it is
 7291 critical that the interface should provide functions to change the scheduling parameters
 7292 dynamically, but, as with many of the other realtime functions, it is important that applications
 7293 use them properly to avoid the possibility of unnecessarily degrading performance.

7294 In providing functions for dynamically changing the scheduling behavior of threads, there were
 7295 two options: provide functions to get and set the individual scheduling parameters of threads,
 7296 or provide a single interface to get and set all the scheduling parameters for a given thread
 7297 simultaneously. Both approaches have merit. Access functions for individual parameters allow
 7298 simpler control of thread scheduling for simple thread scheduling parameters. However, a single
 7299 function for setting all the parameters for a given scheduling policy is required when first setting
 7300 that scheduling policy. Since the single all-encompassing functions are required, it was decided
 7301 to leave the interface as minimal as possible. Note that simpler functions (such as
 7302 *pthread_setprio()* for threads running under the priority-based schedulers) can be easily defined
 7303 in terms of the all-encompassing functions.

7304 If the *pthread_setschedparam()* function executes successfully, it will have set all of the scheduling
 7305 parameter values indicated in *param*; otherwise, none of the scheduling parameters will have
 7306 been modified. This is necessary to ensure that the scheduling of this and all other threads
 7307 continues to be consistent in the presence of an erroneous scheduling parameter.

7308 The [EPERM] error value is included in the list of possible *pthread_setschedparam()* error returns
 7309 as a reflection of the fact that the ability to change scheduling parameters increases risks to the
 7310 implementation and application performance if the scheduling parameters are changed
 7311 improperly. For this reason, and based on some existing practice, it was felt that some
 7312 implementations would probably choose to define specific permissions for changing either a
 7313 thread's own or another thread's scheduling parameters. IEEE Std 1003.1-200x does not include
 7314 portable methods for setting or retrieving permissions, so any such use of permissions is
 7315 completely unspecified.

7316 **Mutex Initialization Scheduling Attributes**

7317 In a priority-driven environment, a direct use of traditional primitives like mutexes and
 7318 condition variables can lead to unbounded priority inversion, where a higher priority thread can
 7319 be blocked by a lower priority thread, or set of threads, for an unbounded duration of time. As a
 7320 result, it becomes impossible to guarantee thread deadlines. Priority inversion can be bounded
 7321 and minimized by the use of priority inheritance protocols. This allows thread deadlines to be
 7322 guaranteed even in the presence of synchronization requirements.

7323 Two useful but simple members of the family of priority inheritance protocols are the basic
 7324 priority inheritance protocol and the priority ceiling protocol emulation. Under the Basic
 7325 Priority Inheritance protocol (governed by the Non-Robust Mutex Priority Inheritance option), a
 7326 thread that is blocking higher priority threads executes at the priority of the highest priority
 7327 thread that it blocks. This simple mechanism allows priority inversion to be bounded by the
 7328 duration of critical sections and makes timing analysis possible.

7329 Under the Priority Ceiling Protocol Emulation protocol (governed by the Thread Priority
 7330 Protection option), each mutex has a priority ceiling, usually defined as the priority of the
 7331 highest priority thread that can lock the mutex. When a thread is executing inside critical
 7332 sections, its priority is unconditionally increased to the highest of the priority ceilings of all the
 7333 mutexes owned by the thread. This protocol has two very desirable properties in uni-processor
 7334 systems. First, a thread can be blocked by a lower priority thread for at most the duration of one
 7335 single critical section. Furthermore, when the protocol is correctly used in a single processor, and
 7336 if threads do not become blocked while owning mutexes, mutual deadlocks are prevented.

7337 The priority ceiling emulation can be extended to multiple processor environments, in which
 7338 case the values of the priority ceilings will be assigned depending on the kind of mutex that is
 7339 being used: local to only one processor, or global, shared by several processors. Local priority
 7340 ceilings will be assigned the usual way, equal to the priority of the highest priority thread that
 7341 may lock that mutex. Global priority ceilings will usually be assigned a priority level higher
 7342 than all the priorities assigned to any of the threads that reside in the involved processors to
 7343 avoid the effect called remote blocking.

7344 **Change the Priority Ceiling of a Mutex**

7345 In order for the priority protect protocol to exhibit its desired properties of bounding priority
 7346 inversion and avoidance of deadlock, it is critical that the ceiling priority of a mutex be the same
 7347 as the priority of the highest thread that can ever hold it, or higher. Thus, if the priorities of the
 7348 threads using such mutexes never change dynamically, there is no need ever to change the
 7349 priority ceiling of a mutex.

7350 However, if a major system mode change results in an altered response time requirement for one
 7351 or more application threads, their priority has to change to reflect it. It will occasionally be the
 7352 case that the priority ceilings of mutexes held also need to change. While changing priority
 7353 ceilings should generally be avoided, it is important that IEEE Std 1003.1-200x provide these
 7354 interfaces for those cases in which it is necessary.

7355 *B.2.9.5 Thread Cancellation*

7356 Many existing threads packages have facilities for canceling an operation or canceling a thread.
 7357 These facilities are used for implementing user requests (such as the CANCEL button in a
 7358 window-based application), for implementing OR parallelism (for example, telling the other
 7359 threads to stop working once one thread has found a forced mate in a parallel chess program), or
 7360 for implementing the ABORT mechanism in Ada.

7361 POSIX programs traditionally have used the signal mechanism combined with either *longjmp()*
 7362 or polling to cancel operations. Many POSIX programmers have trouble using these facilities to
 7363 solve their problems efficiently in a single-threaded process. With the introduction of threads,
 7364 these solutions become even more difficult to use.

7365 The main issues with implementing a cancellation facility are specifying the operation to be
 7366 canceled, cleanly releasing any resources allocated to that operation, controlling when the target
 7367 notices that it has been canceled, and defining the interaction between asynchronous signals and
 7368 cancellation.

7369 **Specifying the Operation to Cancel**

7370 Consider a thread that calls through five distinct levels of program abstraction and then, inside
 7371 the lowest-level abstraction, calls a function that suspends the thread. (An abstraction boundary
 7372 is a layer at which the client of the abstraction sees only the service being provided and can
 7373 remain ignorant of the implementation. Abstractions are often layered, each level of abstraction
 7374 being a client of the lower-level abstraction and implementing a higher-level abstraction.)
 7375 Depending on the semantics of each abstraction, one could imagine wanting to cancel only the
 7376 call that causes suspension, only the bottom two levels, or the operation being done by the entire
 7377 thread. Canceling operations at a finer grain than the entire thread is difficult because threads
 7378 are active and they may be run in parallel on a multi-processor. By the time one thread can make
 7379 a request to cancel an operation, the thread performing the operation may have completed that
 7380 operation and gone on to start another operation whose cancellation is not desired. Thread IDs
 7381 are not reused until the thread has exited, and either it was created with the *Attr detachstate*
 7382 attribute set to *PTHREAD_CREATE_DETACHED* or the *pthread_join()* or *pthread_detach()*
 7383 function has been called for that thread. Consequently, a thread cancellation will never be

7384 misdirected when the thread terminates. For these reasons, the canceling of operations is done at
7385 the granularity of the thread. Threads are designed to be inexpensive enough so that a separate
7386 thread may be created to perform each separately cancelable operation; for example, each
7387 possibly long running user request.

7388 For cancellation to be used in existing code, cancellation scopes and handlers will have to be
7389 established for code that needs to release resources upon cancellation, so that it follows the
7390 programming discipline described in the text.

7391 **A Special Signal Versus a Special Interface**

7392 Two different mechanisms were considered for providing the cancellation interfaces. The first
7393 was to provide an interface to direct signals at a thread and then to define a special signal that
7394 had the required semantics. The other alternative was to use a special interface that delivered the
7395 correct semantics to the target thread.

7396 The solution using signals produced a number of problems. It required the implementation to
7397 provide cancellation in terms of signals whereas a perfectly valid (and possibly more efficient)
7398 implementation could have both layered on a low-level set of primitives. There were so many
7399 exceptions to the special signal (it cannot be used with *kill()*, no POSIX.1 interfaces can be used
7400 with it) that it was clearly not a valid signal. Its semantics on delivery were also completely
7401 different from any existing POSIX.1 signal. As such, a special interface that did not mandate the
7402 implementation and did not confuse the semantics of signals and cancellation was felt to be the
7403 better solution.

7404 **Races Between Cancellation and Resuming Execution**

7405 Due to the nature of cancellation, there is generally no synchronization between the thread
7406 requesting the cancellation of a blocked thread and events that may cause that thread to resume
7407 execution. For this reason, and because excess serialization hurts performance, when both an
7408 event that a thread is waiting for has occurred and a cancellation request has been made and
7409 cancellation is enabled, IEEE Std 1003.1-200x explicitly allows the implementation to choose
7410 between returning from the blocking call or acting on the cancellation request.

7411 **Interaction of Cancellation with Asynchronous Signals**

7412 A typical use of cancellation is to acquire a lock on some resource and to establish a cancellation
7413 cleanup handler for releasing the resource when and if the thread is canceled.

7414 A correct and complete implementation of cancellation in the presence of asynchronous signals
7415 requires considerable care. An implementation has to push a cancellation cleanup handler on the
7416 cancellation cleanup stack while maintaining the integrity of the stack data structure. If an
7417 asynchronously-generated signal is posted to the thread during a stack operation, the signal
7418 handler cannot manipulate the cancellation cleanup stack. As a consequence, asynchronous
7419 signal handlers may not cancel threads or otherwise manipulate the cancellation state of a
7420 thread. Threads may, of course, be canceled by another thread that used a *sigwait()* function to
7421 wait synchronously for an asynchronous signal.

7422 In order for cancellation to function correctly, it is required that asynchronous signal handlers
7423 not change the cancellation state. This requires that some elements of existing practice, such as
7424 using *longjmp()* to exit from an asynchronous signal handler implicitly, be prohibited in cases
7425 where the integrity of the cancellation state of the interrupt thread cannot be ensured.

7426
7427
7428
7429
7430
7431
7432
7433
7434
7435
7436
7437
7438
7439
7440
7441
7442
7443
7444
7445
7446
7447
7448
7449
7450
7451
7452
7453
7454
7455
7456
7457
7458
7459
7460
7461
7462
7463
7464
7465
7466
7467
7468
7469
7470
7471

Thread Cancellation Overview

- Cancelability States

The three possible cancelability states (disabled, deferred, and asynchronous) are encoded into two separate bits ((disable, enable) and (deferred, asynchronous)) to allow them to be changed and restored independently. For instance, short code sequences that will not block sometimes disable cancelability on entry and restore the previous state upon exit. Likewise, long or unbounded code sequences containing no convenient explicit cancellation points will sometimes set the cancelability type to asynchronous on entry and restore the previous value upon exit.

- Cancellation Points

Cancellation points are points inside of certain functions where a thread has to act on any pending cancellation request when cancelability is enabled. For functions in the “shall occur” list, a cancellation check must be performed on every call regardless of whether, absent the cancellation, the call would have blocked. For functions in the “may occur” list, a cancellation check may be performed on some calls but not others; i.e., whether or not a cancellation point occurs when one of these functions is being executed can depend on current conditions.

The idea was considered of allowing implementations to define whether blocking calls such as *read()* should be cancellation points. It was decided that it would adversely affect the design of conforming applications if blocking calls were not cancellation points because threads could be left blocked in an uncancelable state.

There are several important blocking routines that are specifically not made cancellation points:

- *pthread_mutex_lock()*

If *pthread_mutex_lock()* were a cancellation point, every routine that called it would also become a cancellation point (that is, any routine that touched shared state would automatically become a cancellation point). For example, *malloc()*, *free()*, and *rand()* would become cancellation points under this scheme. Having too many cancellation points makes programming very difficult, leading to either much disabling and restoring of cancelability or much difficulty in trying to arrange for reliable cleanup at every possible place.

Since *pthread_mutex_lock()* is not a cancellation point, threads could result in being blocked uninterruptibly for long periods of time if mutexes were used as a general synchronization mechanism. As this is normally not acceptable, mutexes should only be used to protect resources that are held for small fixed lengths of time where not being able to be canceled will not be a problem. Resources that need to be held exclusively for long periods of time should be protected with condition variables.

- *pthread_barrier_wait()*

Canceling a barrier wait will render a barrier unusable. Similar to a barrier timeout (which the standard developers rejected), there is no way to guarantee the consistency of a barrier’s internal data structures if a barrier wait is canceled.

- *pthread_spin_lock()*

As with mutexes, spin locks should only be used to protect resources that are held for small fixed lengths of time where not being cancelable will not be a problem.

Every library routine should specify whether or not it includes any cancellation points. Typically, only those routines that may block or compute indefinitely need to include

7472
7473
7474
7475
7476
7477
7478
7479
7480
7481
7482
7483
7484
7485
7486
7487
7488
7489
7490
7491
7492
7493
7494
7495
7496
7497
7498
7499
7500
7501
7502
7503
7504
7505
7506
7507
7508
7509
7510
7511
7512
7513
7514
7515
7516
7517
7518

cancellation points.

Correctly coded routines only reach cancellation points after having set up a cancellation cleanup handler to restore invariants if the thread is canceled at that point. Being cancelable only at specified cancellation points allows programmers to keep track of actions needed in a cancellation cleanup handler more easily. A thread should only be made asynchronously cancelable when it is not in the process of acquiring or releasing resources or otherwise in a state from which it would be difficult or impossible to recover.

- Thread Cancellation Cleanup Handlers

The cancellation cleanup handlers provide a portable mechanism, easy to implement, for releasing resources and restoring invariants. They are easier to use than signal handlers because they provide a stack of cancellation cleanup handlers rather than a single handler, and because they have an argument that can be used to pass context information to the handler.

The alternative to providing these simple cancellation cleanup handlers (whose only use is for cleaning up when a thread is canceled) is to define a general exception package that could be used for handling and cleaning up after hardware traps and software-detected errors. This was too far removed from the charter of providing threads to handle asynchrony. However, it is an explicit goal of IEEE Std 1003.1-200x to be compatible with existing exception facilities and languages having exceptions.

The interaction of this facility and other procedure-based or language-level exception facilities is unspecified in this version of IEEE Std 1003.1-200x. However, it is intended that it be possible for an implementation to define the relationship between these cancellation cleanup handlers and Ada, C++, or other language-level exception handling facilities.

It was suggested that the cancellation cleanup handlers should also be called when the process exits or calls the *exec* function. This was rejected partly due to the performance problem caused by having to call the cancellation cleanup handlers of every thread before the operation could continue. The other reason was that the only state expected to be cleaned up by the cancellation cleanup handlers would be the intraprocess state. Any handlers that are to clean up the interprocess state would be registered with *atexit()*. There is the orthogonal problem that the *exec* functions do not honor the *atexit()* handlers, but resolving this is beyond the scope of IEEE Std 1003.1-200x.

- Async-Cancel Safety

A function is said to be async-cancel-safe if it is written in such a way that entering the function with asynchronous cancelability enabled will not cause any invariants to be violated, even if a cancellation request is delivered at any arbitrary instruction. Functions that are async-cancel-safe are often written in such a way that they need to acquire no resources for their operation and the visible variables that they may write are strictly limited.

Any routine that gets a resource as a side effect cannot be made async-cancel-safe (for example, *malloc()*). If such a routine were called with asynchronous cancelability enabled, it might acquire the resource successfully, but as it was returning to the client, it could act on a cancellation request. In such a case, the application would have no way of knowing whether the resource was acquired or not.

Indeed, because many interesting routines cannot be made async-cancel-safe, most library routines in general are not async-cancel-safe. Every library routine should specify whether or not it is async-cancel safe so that programmers know which routines can be called from code that is asynchronously cancelable.

7519 IEEE Std 1003.1-2001/Cor 1-2002, item XSH/TC1/D6/8 is applied, adding the *pselect()* function
7520 to the list of functions with cancellation points.

7521 IEEE Std 1003.1-2001/Cor 2-2004, item XSH/TC2/D6/5 is applied, adding the *fdatasync()*
7522 function into the table of functions that shall have cancellation points.

7523 IEEE Std 1003.1-2001/Cor 2-2004, item XSH/TC2/D6/6 is applied, adding the numerous
7524 functions into the table of functions that may have cancellation points.

7525 IEEE Std 1003.1-2001/Cor 2-2004, item XSH/TC2/D6/7 is applied, clarifying the requirements
7526 in Thread Cancellation Cleanup Handlers.

7527 B.2.9.6 Thread Read-Write Locks

7528 Background

7529 Read-write locks are often used to allow parallel access to data on multi-processors, to avoid
7530 context switches on uni-processors when multiple threads access the same data, and to protect
7531 data structures that are frequently accessed (that is, read) but rarely updated (that is, written).
7532 The in-core representation of a file system directory is a good example of such a data structure.
7533 One would like to achieve as much concurrency as possible when searching directories, but limit
7534 concurrent access when adding or deleting files.

7535 Although read-write locks can be implemented with mutexes and condition variables, such
7536 implementations are significantly less efficient than is possible. Therefore, this synchronization
7537 primitive is included in IEEE Std 1003.1-200x for the purpose of allowing more efficient
7538 implementations in multi-processor systems.

7539 Queuing of Waiting Threads

7540 The *pthread_rwlock_unlock()* function description states that one writer or one or more readers
7541 must acquire the lock if it is no longer held by any thread as a result of the call. However, the
7542 function does not specify which thread(s) acquire the lock, unless the Thread Execution
7543 Scheduling option is supported.

7544 The standard developers considered the issue of scheduling with respect to the queuing of
7545 threads blocked on a read-write lock. The question turned out to be whether
7546 IEEE Std 1003.1-200x should require priority scheduling of read-write locks for threads whose
7547 execution scheduling policy is priority-based (for example, SCHED_FIFO or SCHED_RR).
7548 There are tradeoffs between priority scheduling, the amount of concurrency achievable among
7549 readers, and the prevention of writer and/or reader starvation.

7550 For example, suppose one or more readers hold a read-write lock and the following threads
7551 request the lock in the listed order:

```
7552 pthread_rwlock_wrlock() - Low priority thread writer_a
7553 pthread_rwlock_rdlock() - High priority thread reader_a
7554 pthread_rwlock_rdlock() - High priority thread reader_b
7555 pthread_rwlock_rdlock() - High priority thread reader_c
```

7556 When the lock becomes available, should *writer_a* block the high priority readers? Or, suppose a
7557 read-write lock becomes available and the following are queued:

```
7558 pthread_rwlock_rdlock() - Low priority thread reader_a
7559 pthread_rwlock_rdlock() - Low priority thread reader_b
7560 pthread_rwlock_rdlock() - Low priority thread reader_c
7561 pthread_rwlock_wrlock() - Medium priority thread writer_a
7562 pthread_rwlock_rdlock() - High priority thread reader_d
```

7563 If priority scheduling is applied then *reader_d* would acquire the lock and *writer_a* would block

7564 the remaining readers. But should the remaining readers also acquire the lock to increase
 7565 concurrency? The solution adopted takes into account that when the Thread Execution
 7566 Scheduling option is supported, high priority threads may in fact starve low priority threads
 7567 (the application developer is responsible in this case for designing the system in such a way that
 7568 this starvation is avoided). Therefore, IEEE Std 1003.1-200x specifies that high priority readers
 7569 take precedence over lower priority writers. However, to prevent writer starvation from threads
 7570 of the same or lower priority, writers take precedence over readers of the same or lower priority.

7571 Priority inheritance mechanisms are non-trivial in the context of read-write locks. When a high
 7572 priority writer is forced to wait for multiple readers, for example, it is not clear which subset of
 7573 the readers should inherit the writer's priority. Furthermore, the internal data structures that
 7574 record the inheritance must be accessible to all readers, and this implies some sort of
 7575 serialization that could negate any gain in parallelism achieved through the use of multiple
 7576 readers in the first place. Finally, existing practice does not support the use of priority
 7577 inheritance for read-write locks. Therefore, no specification of priority inheritance or priority
 7578 ceiling is attempted. If reliable priority-scheduled synchronization is absolutely required, it can
 7579 always be obtained through the use of mutexes.

7580 **Comparison to *fcntl()* Locks**

7581 The read-write locks and the *fcntl()* locks in IEEE Std 1003.1-200x share a common goal:
 7582 increasing concurrency among readers, thus increasing throughput and decreasing delay.

7583 However, the read-write locks have two features not present in the *fcntl()* locks. First, under
 7584 priority scheduling, read-write locks are granted in priority order. Second, also under priority
 7585 scheduling, writer starvation is prevented by giving writers preference over readers of equal or
 7586 lower priority.

7587 Also, read-write locks can be used in systems lacking a file system, such as those conforming to
 7588 the minimal realtime system profile of IEEE Std 1003.13-1998.

7589 **History of Resolution Issues**

7590 Based upon some balloting objections, early drafts specified the behavior of threads waiting on a
 7591 read-write lock during the execution of a signal handler, as if the thread had not called the lock
 7592 operation. However, this specified behavior would require implementations to establish
 7593 internal signal handlers even though this situation would be rare, or never happen for many
 7594 programs. This would introduce an unacceptable performance hit in comparison to the little
 7595 additional functionality gained. Therefore, the behavior of read-write locks and signals was
 7596 reverted back to its previous mutex-like specification.

7597 *B.2.9.7 Thread Interactions with Regular File Operations*

7598 There is no additional rationale provided for this section.

7599 *B.2.9.8 Use of Application-Managed Thread Stacks*

7600 IEEE Std 1003.1-2001/Cor 2-2004, item XSH/TC2/D6/8 is applied, adding this new section. It
 7601 was added to make it clear that the current standard does not allow an application to determine
 7602 when a stack can be reclaimed. This may be addressed in a future revision.

7603 **B.2.10 Sockets**

7604 The base document for the sockets interfaces in IEEE Std 1003.1-200x is the XNS, Issue 5.2
 7605 specification. This was primarily chosen as it aligns with IPv6. Additional material has been
 7606 added from IEEE Std 1003.1g-2000, notably socket concepts, raw sockets, the *pselect()* function,
 7607 the *socketmark()* function, and the `<sys/select.h>` header.

7608 *B.2.10.1 Address Families*

7609 There is no additional rationale provided for this section.

7610 *B.2.10.2 Addressing*

7611 There is no additional rationale provided for this section.

7612 *B.2.10.3 Protocols*

7613 There is no additional rationale provided for this section.

7614 *B.2.10.4 Routing*

7615 There is no additional rationale provided for this section.

7616 *B.2.10.5 Interfaces*

7617 There is no additional rationale provided for this section.

7618 *B.2.10.6 Socket Types*

7619 The type **socklen_t** was invented to cover the range of implementations seen in the field. The
 7620 intent of **socklen_t** is to be the type for all lengths that are naturally bounded in size; that is, that
 7621 they are the length of a buffer which cannot sensibly become of massive size: network addresses,
 7622 host names, string representations of these, ancillary data, control messages, and socket options
 7623 are examples. Truly boundless sizes are represented by **size_t** as in *read()*, *write()*, and so on.

7624 All **socklen_t** types were originally (in BSD UNIX) of type **int**. During the development of
 7625 IEEE Std 1003.1-200x, it was decided to change all buffer lengths to **size_t**, which appears at face
 7626 value to make sense. When dual mode 32/64-bit systems came along, this choice unnecessarily
 7627 complicated system interfaces because **size_t** (with **long**) was a different size under ILP32 and
 7628 LP64 models. Reverting to **int** would have happened except that some implementations had
 7629 already shipped 64-bit-only interfaces. The compromise was a type which could be defined to be
 7630 any size by the implementation: **socklen_t**.

7631 *B.2.10.7 Socket I/O Mode*

7632 There is no additional rationale provided for this section.

7633 *B.2.10.8 Socket Owner*

7634 There is no additional rationale provided for this section.

7635 *B.2.10.9 Socket Queue Limits*

7636 There is no additional rationale provided for this section.

7637 *B.2.10.10 Pending Error*

7638 There is no additional rationale provided for this section.

7639 *B.2.10.11 Socket Receive Queue*

7640 There is no additional rationale provided for this section.

7641 *B.2.10.12 Socket Out-of-Band Data State*

7642 There is no additional rationale provided for this section.

7643 *B.2.10.13 Connection Indication Queue*

7644 There is no additional rationale provided for this section.

7645 *B.2.10.14 Signals*

7646 There is no additional rationale provided for this section.

7647 *B.2.10.15 Asynchronous Errors*

7648 There is no additional rationale provided for this section.

7649 *B.2.10.16 Use of Options*

7650 There is no additional rationale provided for this section.

7651 *B.2.10.17 Use of Sockets for Local UNIX Connections*

7652 There is no additional rationale provided for this section.

7653 *B.2.10.18 Use of Sockets over Internet Protocols*7654 A raw socket allows privileged users direct access to a protocol; for example, raw access to the
7655 IP and ICMP protocols is possible through raw sockets. Raw sockets are intended for
7656 knowledgeable applications that wish to take advantage of some protocol feature not directly
7657 accessible through the other sockets interfaces.7658 *B.2.10.19 Use of Sockets over Internet Protocols Based on IPv4*

7659 There is no additional rationale provided for this section.

7660 *B.2.10.20 Use of Sockets over Internet Protocols Based on IPv6*7661 The Open Group Base Resolution bwg2001-012 is applied, clarifying that IPv6 implementations
7662 are required to support use of AF_INET6 sockets over IPv4.7663 **B.2.11 Tracing**7664 The organization of the tracing rationale differs from the traditional rationale in that this tracing
7665 rationale text is written against the trace interface as a whole, rather than against the individual
7666 components of the trace interface or the normative section in which those components are
7667 defined. Therefore the sections below do not parallel the sections of normative text in
7668 IEEE Std 1003.1-200x.7669 *B.2.11.1 Objectives*7670 The intended uses of tracing are application-system debugging during system development, as a
7671 “flight recorder” for maintenance of fielded systems, and as a performance measurement tool.
7672 In all of these intended uses, the vendor-supplied computer system and its software are, for this
7673 discussion, assumed error-free; the intent being to debug the user-written and/or third-party
7674 application code, and their interactions. Clearly, problems with the vendor-supplied system and
7675 its software will be uncovered from time to time, but this is a byproduct of the primary activity,
7676 debugging user code.7677 Another need for defining a trace interface in POSIX stems from the objective to provide an
7678 efficient portable way to perform benchmarks. Existing practice shows that such interfaces are
7679 commonly used in a variety of systems but with little commonality. As part of the benchmarking

7680 needs, two aspects within the trace interface must be considered.

7681 The first, and perhaps more important one, is the qualitative aspect.

7682 The second is the quantitative aspect.

7683 • Qualitative Aspect

7684 To better understand this aspect, let us consider an example. Suppose that you want to
 7685 organize a number of actions to be performed during the day. Some of these actions are
 7686 known at the beginning of the day. Some others, which may be more or less important,
 7687 will be triggered by reading your mail. During the day you will make some phone calls
 7688 and synchronously receive some more information. Finally you will receive asynchronous
 7689 phone calls that also will trigger actions. If you, or somebody else, examines your day at
 7690 work, you, or he, can discover that you have not efficiently organized your work. For
 7691 instance, relative to the phone calls you made, would it be preferable to make some of
 7692 these early in the morning? Or to delay some others until the end of the day? Relative to
 7693 the phone calls you have received, you might find that somebody you called in the
 7694 morning has called you 10 times while you were performing some important work. To
 7695 examine, afterwards, your day at work, you record in sequence all the trace events relative
 7696 to your work. This should give you a chance of organizing your next day at work.

7697 This is the qualitative aspect of the trace interface. The user of a system needs to keep a
 7698 trace of particular points the application passes through, so that he can eventually make
 7699 some changes in the application and/or system configuration, to give the application a
 7700 chance of running more efficiently.

7701 • Quantitative Aspect

7702 This aspect concerns primarily realtime applications, where missed deadlines can be
 7703 undesirable. Although there are, in IEEE Std 1003.1-200x, some interfaces useful for such
 7704 applications (timeouts, execution time monitoring, and so on), there are no APIs to aid in
 7705 the tuning of a realtime application's behavior (**timespec** in timeouts, length of message
 7706 queues, duration of driver interrupt service routine, and so on). The tuning of an
 7707 application needs a means of recording timestamped important trace events during
 7708 execution in order to analyze offline, and eventually, to tune some realtime features
 7709 (redesign the system with less functionalities, readjust timeouts, redesign driver interrupts,
 7710 and so on).

7711 **Detailed Objectives**

7712 Objectives were defined to build the trace interface and are kept for historical interest. Although
 7713 some objectives are not fully respected in this trace interface, the concept of the POSIX trace
 7714 interface assumes the following points:

- 7715 1. It must be possible to trace both system and user trace events concurrently.
- 7716 2. It must be possible to trace per-process trace events and also to trace system trace events
 7717 which are unrelated to any particular process. A per-process trace event is either user-
 7718 initiated or system-initiated.
- 7719 3. It must be possible to control tracing on a per-process basis from either inside or outside
 7720 the process.
- 7721 4. It must be possible to control tracing on a per-thread basis from inside the enclosing
 7722 process.

- 7723 5. Trace points must be controllable by trace event type ID from inside and outside of the
7724 process. Multiple trace points can have the same trace event type ID, and will be
7725 controlled jointly.
- 7726 6. Recording of trace events is dependent on both trace event type ID and the
7727 process/thread. Both must be enabled in order to record trace events. System trace events
7728 may or may not be handled differently.
- 7729 7. The API must not mandate the ability to control tracing for more than one process at the
7730 same time.
- 7731 8. There is no objective for trace control on anything bigger than a process; for example,
7732 group or session.
- 7733 9. Trace propagation and control:
- 7734 a. Trace propagation across *fork()* is optional; the default is to not trace a child
7735 process.
- 7736 b. Trace control must span *pthread_create()* operations; that is, if a process is being
7737 traced, any thread will be traced as well if this thread allows tracing. The default is
7738 to allow tracing.
- 7739 10. Trace control must not span *exec* or *posix_spawn()* operations.
- 7740 11. A triggering API is not required. The triggering API is the ability to command or stop
7741 tracing based on the occurrence of a specific trace event other than a
7742 POSIX_TRACE_START trace event or a POSIX_TRACE_STOP trace event.
- 7743 12. Trace log entries must have timestamps of implementation-defined resolution.
7744 Implementations are exhorted to support at least microsecond resolution. When a trace
7745 log entry is retrieved, it must have timestamp, PC address, PID, and TID of the entity that
7746 generated the trace event.
- 7747 13. Independently developed code should be able to use trace facilities without coordination
7748 and without conflict.
- 7749 14. Even if the trace points in the trace calls are not unique, the trace log entries (after any
7750 processing) must be uniquely identified as to trace point.
- 7751 15. There must be a standard API to read the trace stream.
- 7752 16. The format of the trace stream and the trace log is opaque and unspecified.
- 7753 17. It must be possible to read a completed trace, if recorded on some suitable non-volatile
7754 storage, even subsequent to a power cycle or subsequent cold boot of the system.
- 7755 18. Support of analysis of a trace log while it is being formed is implementation-defined.
- 7756 19. The API must allow the application to write trace stream identification information into
7757 the trace stream and to be able to retrieve it, without it being overwritten by trace entries,
7758 even if the trace stream is full.
- 7759 20. It must be possible to specify the destination of trace data produced by trace events.
- 7760 21. It must be possible to have different trace streams, and for the tracing enabled by one
7761 trace stream to be completely independent of the tracing of another trace stream.
- 7762 22. It must be possible to trace events from threads in different CPUs.
- 7763 23. The API must support one or more trace streams per-system, and one or more trace
7764 streams per-process, up to an implementation-defined set of per-system and per-process
7765 maximums.

- 7766 24. It must be possible to determine the order in which the trace events happened, without
7767 necessarily depending on the clock, up to an implementation-defined time resolution.
- 7768 25. For performance reasons, the trace event point call(s) must be implementable as a macro
7769 (see the ISO POSIX-1: 1996 standard, 1.3.4, Statement 2).
- 7770 26. IEEE Std 1003.1-200x must not define the trace points which a conforming system must
7771 implement, except for trace points used in the control of tracing.
- 7772 27. The APIs must be thread-safe, and trace points should be lock-free (that is, not require a
7773 lock to gain exclusive access to some resource).
- 7774 28. The user-provided information associated with a trace event is variable-sized, up to some
7775 maximum size.
- 7776 29. Bounds on record and trace stream sizes:
- 7777 a. The API must permit the application to declare the upper bounds on the length of
7778 an application data record. The system must return the limit it used. The limit used
7779 may be smaller than requested.
- 7780 b. The API must permit the application to declare the upper bounds on the size of
7781 trace streams. The system must return the limit it used. The limit used may be
7782 different, either larger or smaller, than requested.
- 7783 30. The API must be able to pass any fundamental data type, and a structured data type
7784 composed only of fundamental types. The API must be able to pass data by reference,
7785 given only as an address and a length. Fundamental types are the POSIX.1 types (see the
7786 `<sys/types.h>` header) plus those defined in the ISO C standard.
- 7787 31. The API must apply the POSIX notions of ownership and permission to recorded trace
7788 data, corresponding to the sources of that data.

7789 Comments on Objectives

7790 **Note:** In the following comments, numbers in square brackets refer to the above objectives.

7791 It is necessary to be able to obtain a trace stream for a complete activity. Thus there is a
7792 requirement to be able to trace both application and system trace events. A per-process trace
7793 event is either user-initiated, like the `write()` function, or system-initiated, like a timer expiration.
7794 There is also a need to be able to trace the activity of an entire process even when it has threads
7795 in multiple CPUs. To avoid excess trace activity, it is necessary to be able to control tracing on a
7796 trace event type basis.
7797 [Objectives 1,2,5,22]

7798 There is a need to be able to control tracing on a per-process basis, both from inside and outside
7799 the process; that is, a process can start a trace activity on itself or any other process. There is also
7800 the perceived need to allow the definition of a maximum number of trace streams per system.
7801 [Objectives 3,23]

7802 From within a process, it is necessary to be able to control tracing on a per-thread basis. This
7803 provides an additional filtering capability to keep the amount of traced data to a minimum. It
7804 also allows for less ambiguity as to the origin of trace events. It is recognized that thread-level
7805 control is only valid from within the process itself. It is also desirable to know the maximum
7806 number of trace streams per process that can be started. The API should not require thread
7807 synchronization or mandate priority inversions that would cause the thread to block. However,
7808 the API must be thread-safe.
7809 [Objectives 4,23,24,27]

7810 There was no perceived objective to control tracing on anything larger than a process; for

7811 example, a group or session. Also, the ability to start or stop a trace activity on multiple
7812 processes atomically may be very difficult or cumbersome in some implementations.
7813 [Objectives 6,8]

7814 It is also necessary to be able to control tracing by trace event type identifier, sometimes called a
7815 trace hook ID. However, there is no mandated set of system trace events, since such trace points
7816 are implementation-defined. The API must not require from the operating system facilities that
7817 are not standard.
7818 [Objectives 6,26]

7819 Trace control must span *fork()* and *pthread_create()*. If not, there will be no way to ensure that an
7820 application's activity is entirely traced. The newly forked child would not be able to turn on its
7821 tracing until after it obtained control after the fork, and trace control externally would be even
7822 more problematic.
7823 [Objective 9]

7824 Since *exec* and *posix_spawn()* represent a complete change in the execution of a task (a new
7825 program), trace control need not persist over an *exec* or *posix_spawn()*.
7826 [Objective 10]

7827 Where trace activities are started on multiple processes, these trace activities should not interfere
7828 with each other.
7829 [Objective 21]

7830 There is no need for a triggering objective, primarily for performance reasons; see also [Section](#)
7831 [B.2.11.8](#) (on page 201), rationale on triggering.
7832 [Objective 11]

7833 It must be possible to determine the origin of each traced event. The process and thread
7834 identifiers for each trace event are needed. Also there was a perceived need for a user-specifiable
7835 origin, but it was felt that this would create too much overhead.
7836 [Objectives 12,14]

7837 An allowance must be made for trace points to come embedded in software components from
7838 several different sources and vendors without requiring coordination.
7839 [Objective 13]

7840 There is a requirement to be able to uniquely identify trace points that may have the same trace
7841 stream identifier. This is only necessary when a trace report is produced.
7842 [Objectives 12,14]

7843 Tracing is a very performance-sensitive activity, and will therefore likely be implemented at a
7844 low level within the system. Hence the interface must not mandate any particular buffering or
7845 storage method. Therefore, a standard API is needed to read a trace stream. Also the interface
7846 must not mandate the format of the trace data, and the interface must not assume a trace storage
7847 method. Due to the possibility of a monolithic kernel and the possible presence of multiple
7848 processes capable of running trace activities, the two kinds of trace events may be stored in two
7849 separate streams for performance reasons. A mandatory dump mechanism, common in some
7850 existing practice, has been avoided to allow the implementation of this set of functions on small
7851 realtime profiles for which the concept of a file system is not defined. The trace API calls should
7852 be implemented as macros.
7853 [Objectives 15,16,25,30]

7854 Since a trace facility is a valuable service tool, the output (or log) of a completed trace stream
7855 that is written to permanent storage must be readable on other systems of the type that
7856 produced the trace log. Note that there is no objective to be able to interpret a trace log that was
7857 not successfully completed.
7858 [Objectives 17,18,19]

7859 For trace streams written to permanent storage, a way to specify the destination of the trace
7860 stream is needed.
7861 [Objective 20]

7862 There is a requirement to be able to depend on the ordering of trace events up to some
7863 implementation-defined time interval. For example, there is a need to know the time period
7864 during which, if trace events are closer together, their ordering is unspecified. Events that occur
7865 within an interval smaller than this resolution may or may not be read back in the correct order.
7866 [Objective 24]

7867 The application should be able to know how much data can be traced. When trace event types
7868 can be filtered, the application should be able to specify the approximate maximum amount of
7869 data that will be traced in a trace event so resources can be more efficiently allocated.
7870 [Objectives 28,29]

7871 Users should not be able to trace data to which they would not normally have access. System
7872 trace events corresponding to a process/thread should be associated with the ownership of that
7873 process/thread.
7874 [Objective 31]

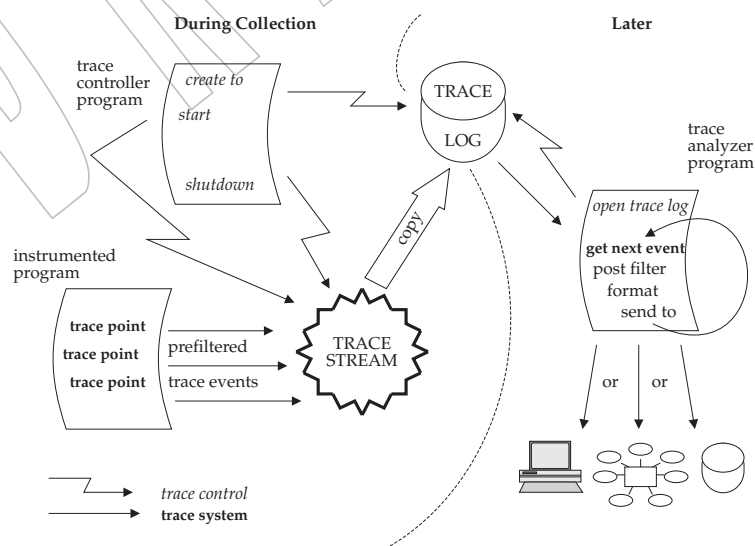
7875 B.2.11.2 Trace Model

7876 Introduction

7877 The model is based on two base entities: the "Trace Stream" and the "Trace Log", and a recorded
7878 unit called the "Trace Event". The possibility of using Trace Streams and Trace Logs separately
7879 gives two use dimensions and solves both the performance issue and the full-information
7880 system issue. In the case of a trace stream without log, specific information, although reduced in
7881 quantity, is required to be registered, in a possibly small realtime system, with as little overhead
7882 as possible. The Trace Log option has been added for small realtime systems. In the case of a
7883 trace stream with log, considerable complex application-specific information needs to be
7884 collected.

7885 Trace Model Description

7886 The trace model can be examined for three different subfunctions: Application Instrumentation,
7887 Trace Operation Control, and Trace Analysis.



7888

Figure B-2 Trace System Overview: for Offline Analysis

7889

Each of these subfunctions requires specific characteristics of the trace mechanism API.

7890

- Application Instrumentation

7891

When instrumenting an application, the programmer is not concerned about the future use of the trace events in the trace stream or the trace log, the full policy of the trace stream, or the eventual pre-filtering of trace events. But he is concerned about the correct determination of the specific trace event type identifier, regardless of how many independent libraries are used in the same user application; see [Figure B-2](#) and [Figure B-3](#) (on page 185).

7892

7893

7894

7895

7896

7897

This trace API provides the necessary operations to accomplish this subfunction. This is done by providing functions to associate a programmer-defined name with an implementation-defined trace event type identifier (see the `posix_trace_eventid_open()` function), and to send this trace event into a potential trace stream (see the `posix_trace_event()` function).

7898

7899

7900

7901

7902

- Trace Operation Control

7903

When controlling the recording of trace events in a trace stream, the programmer is concerned with the correct initialization of the trace mechanism (that is, the sizing of the trace stream), the correct retention of trace events in a permanent storage, the correct dynamic recording of trace events, and so on.

7904

7905

7906

7907

This trace API provides the necessary material to permit this efficiently. This is done by providing functions to initialize a new trace stream, and optionally a trace log:

7908

7909

- Trace Stream Attributes Object Initialization (see `posix_trace_attr_init()`)

7910

- Functions to Retrieve or Set Information About a Trace Stream (see `posix_trace_attr_getgenversion()`)

7911

7912

- Functions to Retrieve or Set the Behavior of a Trace Stream (see `posix_trace_attr_getinherited()`)

7913

7914

- Functions to Retrieve or Set Trace Stream Size Attributes (see `posix_trace_attr_getmaxusereventsize()`)

7915

7916

- Trace Stream Initialization, Flush, and Shutdown from a Process (see `posix_trace_create()`)

7917

7918

- Clear Trace Stream and Trace Log (see `posix_trace_clear()`)

7919

To select the trace event types that are to be traced:

7920

- Manipulate Trace Event Type Identifier (see `posix_trace_trid_eventid_open()`)

7921

- Iterate over a Mapping of Trace Event Type (see `posix_trace_eventtypelist_getnext_id()`)

7922

- Manipulate Trace Event Type Sets (see `posix_trace_eventset_empty()`)

7923

- Set Filter of an Initialized Trace Stream (see `posix_trace_set_filter()`)

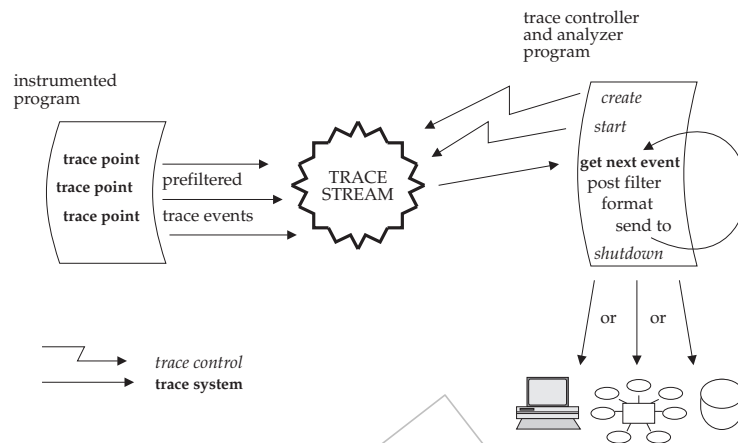
7924

To control the execution of an active trace stream:

7925

- Trace Start and Stop (see `posix_trace_start()`)

7926 — Functions to Retrieve the Trace Attributes or Trace Statuses (see
7927 *posix_trace_get_attr()*)



7928 **Figure B-3** Trace System Overview: for Online Analysis

7929 • Trace Analysis

7930 Once correctly recorded, on permanent storage or not, an ultimate activity consists of the
7931 analysis of the recorded information. If the recorded data is on permanent storage, a
7932 specific open operation is required to associate a trace stream to a trace log.

7933 The first intent of the group was to request the presence of a system identification structure
7934 in the trace stream attribute. This was, for the application, to allow some portable way to
7935 process the recorded information. However, there is no requirement that the **utsname**
7936 structure, on which this system identification was based, be portable from one machine to
7937 another, so the contents of the attribute cannot be interpreted correctly by an application
7938 conforming to IEEE Std 1003.1-200x.

7939 This modification has been incorporated and requests that some unspecified information
7940 be recorded in the trace log in order to fail opening it if the analysis process and the
7941 controller process were running in different types of machine, but does not request that
7942 this information be accessible to the application. This modification has implied a
7943 modification in the *posix_trace_open()* function error code returns.

7944 This trace API provides functions to:

- 7945 — Extract trace stream identification attributes (see *posix_trace_attr_getgenversion()*)
- 7946 — Extract trace stream behavior attributes (see *posix_trace_attr_getinherited()*)
- 7947 — Extract trace event, stream, and log size attributes (see
7948 *posix_trace_attr_getmaxusereventsizesize()*)
- 7949 — Look up trace event type names (see *posix_trace_eventid_get_name()*)
- 7950 — Iterate over trace event type identifiers (see *posix_trace_eventtypelist_getnext_id()*)
- 7951 — Open, rewind, and close a trace log (see *posix_trace_open()*)
- 7952 — Read trace stream attributes and status (see *posix_trace_get_attr()*)

— Read trace events (see *posix_trace_getnext_event()*)

Due to the following two reasons:

1. The requirement that the trace system must not add unacceptable overhead to the traced process and so that the trace event point execution must be fast
2. The traced application does not care about tracing errors

the trace system cannot return any internal error to the application. Internal error conditions can range from unrecoverable errors that will force the active trace stream to abort, to small errors that can affect the quality of tracing without aborting the trace stream. The group decided to define a system trace event to report to the analysis process such internal errors. It is not the intention of IEEE Std 1003.1-200x to require an implementation to report an internal error that corrupts or terminates tracing operation. The implementor is free to decide which internal documented errors, if any, the trace system is able to report.

States of a Trace Stream

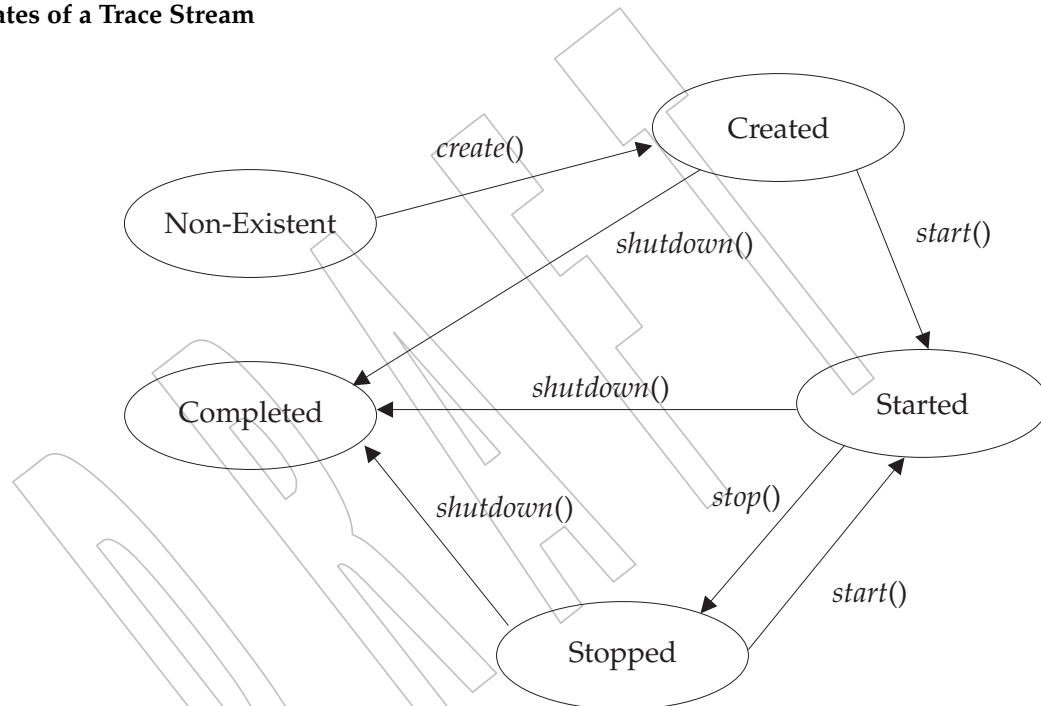


Figure B-4 Trace System Overview: States of a Trace Stream

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.

After a call to *posix_trace_shutdown()*, the trace stream is in the state COMPLETED. The trace

7977 stream no longer exists but, if the Trace Log option is supported, all the information contained in
 7978 it has been logged. If a log object has not been associated with the trace stream at the creation, it
 7979 is the responsibility of the trace controller process to not shut the trace stream down while trace
 7980 events remain to be read in the stream.

7981 **Tracing All Processes**

7982 Some implementations have a tracing subsystem with the ability to trace all processes. This is
 7983 useful to debug some types of device drivers such as those for ATM or X25 adapters. These
 7984 types of adapters are used by several independent processes, that are not issued from the same
 7985 process.

7986 The POSIX trace interface does not define any constant or option to create a trace stream tracing
 7987 all processes. POSIX.1 does not prevent this type of implementation and an implementor is free
 7988 to add this capability. Nevertheless, the trace interface allows tracing of all the system trace
 7989 events and all the processes issued from the same process.

7990 If such a tracing system capability has to be implemented, when a trace stream is created, it is
 7991 recommended that a constant named `POSIX_TRACE_ALLPROC` be used instead of the process
 7992 identifier in the argument of the `posix_trace_create()` or `posix_trace_create_withlog()` function. A
 7993 possible value for `POSIX_TRACE_ALLPROC` may be `-1` instead of a real process identifier.

7994 The implementor has to be aware that there is some impact on the tracing behavior as defined in
 7995 the POSIX trace interface. For example:

- 7996 • If the default value for the inheritance attribute is set to
 7997 `POSIX_TRACE_CLOSE_FOR_CHILD`, the implementation has to stop tracing for the child
 7998 process.
- 7999 • The trace controller which is creating this type of trace stream must have the appropriate
 8000 privilege to trace all the processes.

8001 **Trace Storage**

8002 The model is based on two types of trace events: system trace events and user-defined trace
 8003 events. The internal representation of trace events is implementation-defined, and so the
 8004 implementor is free to choose the more suitable, practical, and efficient way to design the
 8005 internal management of trace events. For the timestamping operation, the model does not
 8006 impose the `CLOCK_REALTIME` or any other clock. The buffering allocation and operation
 8007 follow the same principle. The implementor is free to use one or more buffers to record trace
 8008 events; the interface assumes only a logical trace stream of sequentially recorded trace events.
 8009 Regarding flushing of trace events, the interface allows the definition of a trace log object which
 8010 typically can be a file. But the group was also aware of defining functions to permit the use of
 8011 this interface in small realtime systems, which may not have general file system capabilities. For
 8012 instance, the three functions `posix_trace_getnext_event()` (blocking),
 8013 `posix_trace_timedgetnext_event()` (blocking with timeout), and `posix_trace_trygetnext_event()` (non-
 8014 blocking) are proposed to read the recorded trace events.

8015 The policy to be used when the trace stream becomes full also relies on common practice:

- 8016 • For an active trace stream, the `POSIX_TRACE_LOOP` trace stream policy permits
 8017 automatic overrun (overwrite of oldest trace events) while waiting for some user-defined
 8018 condition to cause tracing to stop. By contrast, the `POSIX_TRACE_UNTIL_FULL` trace
 8019 stream policy requires the system to stop tracing when the trace stream is full. However, if
 8020 the trace stream that is full is at least partially emptied by a call to the `posix_trace_flush()`
 8021 function or by calls to the `posix_trace_getnext_event()` function, the trace system will
 8022 automatically resume tracing.

8023 If the Trace Log option is supported, the operation of the POSIX_TRACE_FLUSH policy is
 8024 an extension of the POSIX_TRACE_UNTIL_FULL policy. The automatic free operation (by
 8025 flushing to the associated trace log) is added.

8026 • If a log is associated with the trace stream and this log is a regular file, these policies also
 8027 apply for the log. One more policy, POSIX_TRACE_APPEND, is defined to allow
 8028 indefinite extension of the log. Since the log destination can be any device or pseudo-
 8029 device, the implementation may not be able to manipulate the destination as required by
 8030 IEEE Std 1003.1-200x. For this reason, the behavior of the log full policy may be
 8031 unspecified depending on the trace log type.

8032 The current trace interface does not define a service to preallocate space for a trace log file,
 8033 because this space can be preallocated by means of a call to the *posix_fallocate()* function.
 8034 This function could be called after the file has been opened, but before the trace stream is
 8035 created. The *posix_fallocate()* function ensures that any required storage for regular file data
 8036 is allocated on the file system storage media. If *posix_fallocate()* returns successfully,
 8037 subsequent writes to the specified file data will not fail due to the lack of free space on the
 8038 file system storage media. Besides trace events, a trace stream also includes trace attributes
 8039 and the mapping from trace event names to trace event type identifiers. The implementor
 8040 is free to choose how to store the trace attributes and the trace event type map, but must
 8041 ensure that this information is not lost when a trace stream overrun occurs.

8042 B.2.11.3 Trace Programming Examples

8043 Several programming examples are presented to show the code of the different possible
 8044 subfunctions using a trace subsystem. All these programs need to include the **<trace.h>** header.
 8045 In the examples shown, error checking is omitted for more simplicity.

8046 Trace Operation Control

8047 These examples show the creation of a trace stream for another process; one which is already
 8048 trace instrumented. All the default trace stream attributes are used to simplify programming in
 8049 the first example. The second example shows more possibilities.

8050 First Example

```
8051 /* Caution. Error checks omitted */
8052 {
8053     trace_attr_t attr;
8054     pid_t pid = traced_process_pid;
8055     int fd;
8056     trace_id_t trid;
8057     - - - - -
8058     /* Initialize trace stream attributes */
8059     posix_trace_attr_init(&attr);
8060     /* Open a trace log */
8061     fd=open("/tmp/mytracelog",...);
8062     /*
8063      * Create a new trace associated with a log
8064      * and with default attributes
8065      */
8066     posix_trace_create_withlog(pid, &attr, fd, &trid);
8067     /* Trace attribute structure can now be destroyed */
8068     posix_trace_attr_destroy(&attr);
```

```

8069         /* Start of trace event recording */
8070         posix_trace_start(trid);
8071         - - - - -
8072         - - - - -
8073         /* Duration of tracing */
8074         - - - - -
8075         - - - - -
8076         /* Stop and shutdown of trace activity */
8077         posix_trace_shutdown(trid);
8078         - - - - -
8079     }

```

8080

Second Example

8081 Between the initialization of the trace stream attributes and the creation of the trace stream, these
8082 trace stream attributes may be modified; see [Trace Stream Attribute Manipulation](#) for a specific
8083 programming example. Between the creation and the start of the trace stream, the event filter
8084 may be set; after the trace stream is started, the event filter may be changed. The setting of an
8085 event set and the change of a filter is shown in [Create a Trace Event Type Set and Change the
8086 Trace Event Type Filter](#) (on page 193).

```

8087     /* Caution. Error checks omitted */
8088     {
8089         trace_attr_t attr;
8090         pid_t pid = traced_process_pid;
8091         int fd;
8092         trace_id_t trid;
8093         - - - - -
8094         /* Initialize trace stream attributes */
8095         posix_trace_attr_init(&attr);
8096         /* Attr default may be changed at this place; see example */
8097         - - - - -
8098         /* Create and open a trace log with R/W user access */
8099         fd=open("/tmp/mytracelog",O_WRONLY|O_CREAT,S_IRUSR|S_IWUSR);
8100         /* Create a new trace associated with a log */
8101         posix_trace_create_withlog(pid, &attr, fd, &trid);
8102         /*
8103          * If the Trace Filter option is supported
8104          * trace event type filter default may be changed at this place;
8105          * see example about changing the trace event type filter
8106          */
8107         posix_trace_start(trid);
8108         - - - - -
8109
8110         /*
8111          * If you have an uninteresting part of the application
8112          * you can stop temporarily.
8113          *
8114          * posix_trace_stop(trid);
8115          * - - - - -
8116          * - - - - -
8117          * posix_trace_start(trid);
8118          */
8119         - - - - -

```

```

8119     /*
8120     * If the Trace Filter option is supported
8121     * the current trace event type filter can be changed
8122     * at any time (see example about how to set
8123     * a trace event type filter)
8124     */
8125     - - - - -
8126     /* Stop the recording of trace events */
8127     posix_trace_stop(trid);
8128     /* Shutdown the trace stream */
8129     posix_trace_shutdown(trid);
8130     /*
8131     * Destroy trace stream attributes; attr structure may have
8132     * been used during tracing to fetch the attributes
8133     */
8134     posix_trace_attr_destroy(&attr);
8135     - - - - -
8136 }

```

8137 Application Instrumentation

8138 This example shows an instrumented application. The code is included in a block of instructions,
8139 perhaps a function from a library. Possibly in an initialization part of the instrumented
8140 application, two user trace event names are mapped to two trace event type identifiers
8141 (function *posix_trace_eventid_open()*). Then two trace points are programmed.

```

8142 /* Caution. Error checks omitted */
8143 {
8144     trace_event_id_t eventid1, eventid2;
8145     - - - - -
8146     /* Initialization of two trace event type ids */
8147     posix_trace_eventid_open("my_first_event",&eventid1);
8148     posix_trace_eventid_open("my_second_event",&eventid2);
8149     - - - - -
8150     - - - - -
8151     - - - - -
8152     /* Trace point */
8153     posix_trace_event(eventid1,NULL,0);
8154     - - - - -
8155     /* Trace point */
8156     posix_trace_event(eventid2,NULL,0);
8157     - - - - -
8158 }

```


8159 **Trace Analyzer**

8160 This example shows the manipulation of a trace log resulting from the dumping of a completed
 8161 trace stream. All the default attributes are used to simplify programming, and data associated
 8162 with a trace event is not shown in the first example. The second example shows more
 8163 possibilities.

8164 **First Example**

```

8165 /* Caution. Error checks omitted */
8166 {
8167     int fd;
8168     trace_id_t trid;
8169     posix_trace_event_info trace_event;
8170     char trace_event_name[TRACE_EVENT_NAME_MAX];
8171     int return_value;
8172     size_t returndatasize;
8173     int lost_event_number;
8174
8175     - - - - -
8176     /* Open an existing trace log */
8177     fd=open("/tmp/tracelog", O_RDONLY);
8178     /* Open a trace stream on the open log */
8179     posix_trace_open(fd, &trid);
8180     /* Read a trace event */
8181     posix_trace_getnext_event(trid, &trace_event,
8182         NULL, 0, &returndatasize,&return_value);
8183     /* Read and print all trace event names out in a loop */
8184     while (return_value == NULL)
8185     {
8186         /*
8187          * Get the name of the trace event associated
8188          * with trid trace ID
8189          */
8190         posix_trace_eventid_get_name(trid, trace_event.event_id,
8191             trace_event_name);
8192         /* Print the trace event name out */
8193         printf("%s\n",trace_event_name);
8194         /* Read a trace event */
8195         posix_trace_getnext_event(trid, &trace_event,
8196             NULL, 0, &returndatasize,&return_value);
8197     }
8198     /* Close the trace stream */
8199     posix_trace_close(trid);
8200     /* Close the trace log */
8201     close(fd);
8202 }
  
```

8202 **Second Example**

8203 The complete example includes the two other examples in [Retrieve Information from a Trace Log](#) and in [Retrieve the List of Trace Event Types Used in a Trace Log](#) (on page 195). For
 8204 example, the *maxdatasize* variable is set in [Retrieve the List of Trace Event Types Used in a Trace Log](#)
 8205 (on page 195).
 8206

```

8207 /* Caution. Error checks omitted */
8208 {
8209     int fd;
8210     trace_id_t trid;
8211     posix_trace_event_info trace_event;
8212     char trace_event_name[TRACE_EVENT_NAME_MAX];
8213     char * data;
8214     size_t maxdatasize=1024, returndatasize;
8215     int return_value;
8216     - - - - -
8217     /* Open an existing trace log */
8218     fd=open("/tmp/tracelog", O_RDONLY);
8219     /* Open a trace stream on the open log */
8220     posix_trace_open( fd, &trid);
8221     /*
8222      * Retrieve information about the trace stream which
8223      * was dumped in this trace log (see example)
8224      */
8225     - - - - -
8226     /* Allocate a buffer for trace event data */
8227     data=(char *)malloc(maxdatasize);
8228     /*
8229      * Retrieve the list of trace events used in this
8230      * trace log (see example)
8231      */
8232     - - - - -
8233     /* Read and print all trace event names and data out in a loop */
8234     while (1)
8235     {
8236         posix_trace_getnext_event(trid, &trace_event,
8237             data, maxdatasize, &returndatasize,&return_value);
8238         if (return_value != NULL) break;
8239         /*
8240          * Get the name of the trace event type associated
8241          * with trid trace ID
8242          */
8243         posix_trace_eventid_get_name(trid, trace_event.event_id,
8244             trace_event_name);
8245         {
8246             int i;
8247             /* Print the trace event name out */
8248             printf("%s: ", trace_event_name);
8249             /* Print the trace event data out */
8250             for (i=0; i<returndatasize, i++) printf("%02.2X",
8251                 (unsigned char)data[i]);

```

```

8252         printf("\n");
8253     }
8254 }
8255     /* Close the trace stream */
8256     posix_trace_close(trid);
8257     /* The buffer data is deallocated */
8258     free(data);
8259     /* Now the file can be closed */
8260     close(fd);
8261 }

```

8262 Several Programming Manipulations

8263 The following examples show some typical sets of operations needed in some contexts.

8264 Trace Stream Attribute Manipulation

8265 This example shows the manipulation of a trace stream attribute object in order to change the
8266 default value provided by a previous `posix_trace_attr_init()` call.

```

8267     /* Caution. Error checks omitted */
8268     {
8269         trace_attr_t attr;
8270         size_t logsize=100000;
8271         - - - - -
8272         /* Initialize trace stream attributes */
8273         posix_trace_attr_init(&attr);
8274         /* Set the trace name in the attributes structure */
8275         posix_trace_attr_setname(&attr, "my_trace");
8276         /* Set the trace full policy */
8277         posix_trace_attr_setstreamfullpolicy(&attr, POSIX_TRACE_LOOP);
8278         /* Set the trace log size */
8279         posix_trace_attr_setlogsize(&attr, logsize);
8280         - - - - -
8281     }

```

8282 Create a Trace Event Type Set and Change the Trace Event Type Filter

8283 This example is valid only if the Trace Event Filter option is supported. This example shows the
8284 manipulation of a trace event type set in order to change the trace event type filter for an
8285 existing active trace stream, which may be just-created, running, or suspended. Some sets of
8286 trace event types are well-known, such as the set of trace event types not associated with a
8287 process, some trace event types are just-built trace event types for this trace stream; one trace
8288 event type is the predefined trace event error type which is deleted from the trace event type set.

```

8289     /* Caution. Error checks omitted */
8290     {
8291         trace_id_t trid = existing_trace;
8292         trace_event_set_t set;
8293         trace_event_id_t trace_event1, trace_event2;
8294         - - - - -
8295         /* Initialize to an empty set of trace event types */
8296         /* (not strictly required because posix_trace_event_set_fill() */
8297         /* will ignore the prior contents of the event set.) */

```

```

8298     posix_trace_eventset_emptyset(&set);
8299     /*
8300     * Fill the set with all system trace events
8301     * not associated with a process
8302     */
8303     posix_trace_eventset_fill(&set, POSIX_TRACE_WOPID_EVENTS);
8304     /*
8305     * Get the trace event type identifier of the known trace event name
8306     * my_first_event for the trid trace stream
8307     */
8308     posix_trace_trid_eventid_open(trid, "my_first_event", &trace_event1);
8309     /* Add the set with this trace event type identifier */
8310     posix_trace_eventset_add_event(trace_event1, &set);
8311     /*
8312     * Get the trace event type identifier of the known trace event name
8313     * my_second_event for the trid trace stream
8314     */
8315     posix_trace_trid_eventid_open(trid, "my_second_event", &trace_event2);
8316     /* Add the set with this trace event type identifier */
8317     posix_trace_eventset_add_event(trace_event2, &set);
8318     - - - - -
8319     /* Delete the system trace event POSIX_TRACE_ERROR from the set */
8320     posix_trace_eventset_del_event(POSIX_TRACE_ERROR, &set);
8321     - - - - -
8322     /* Modify the trace stream filter making it equal to the new set */
8323     posix_trace_set_filter(trid, &set, POSIX_TRACE_SET_EVENTSET);
8324     - - - - -
8325     /*
8326     * Now trace_event1, trace_event2, and all system trace event types
8327     * not associated with a process, except for the POSIX_TRACE_ERROR
8328     * system trace event type, are filtered out of (not recorded in) the
8329     * existing trace stream.
8330     */
8331 }

```

8332 Retrieve Information from a Trace Log

8333 This example shows how to extract information from a trace log, the dump of a trace stream.
8334 This code:

- 8335 • Asks if the trace stream has lost trace events
- 8336 • Extracts the information about the version of the trace subsystem which generated this
8337 trace log
- 8338 • Retrieves the maximum size of trace event data; this may be used to dynamically allocate
8339 an array for extracting trace event data from the trace log without overflow

```

8340 /* Caution. Error checks omitted */
8341 {
8342     struct posix_trace_status_info statusinfo;
8343     trace_attr_t attr;
8344     trace_id_t trid = existing_trace;

```

```

8345     size_t maxdatasize;
8346     char genversion[TRACE_NAME_MAX];
8347     - - - - -
8348     /* Get the trace stream status */
8349     posix_trace_get_status(trid, &statusinfo);
8350     /* Detect an overrun condition */
8351     if (statusinfo.posix_stream_overrun_status == POSIX_TRACE_OVERRUN)
8352         printf("trace events have been lost\n");
8353
8354     /* Get attributes from the trid trace stream */
8355     posix_trace_get_attr(trid, &attr);
8356     /* Get the trace generation version from the attributes */
8357     posix_trace_attr_getgenversion(&attr, genversion);
8358     /* Print the trace generation version out */
8359     printf("Information about Trace Generator:%s\n",genversion);
8360
8361     /* Get the trace event max data size from the attributes */
8362     posix_trace_attr_getmaxdatasize(&attr, &maxdatasize);
8363     /* Print the trace event max data size out */
8364     printf("Maximum size of associated data:%d\n",maxdatasize);
8365     /* Destroy the trace stream attributes */
8366     posix_trace_attr_destroy(&attr);
8367 }

```

Retrieve the List of Trace Event Types Used in a Trace Log

This example shows the retrieval of a trace stream's trace event type list. This operation may be very useful if you are interested only in tracking the type of trace events in a trace log.

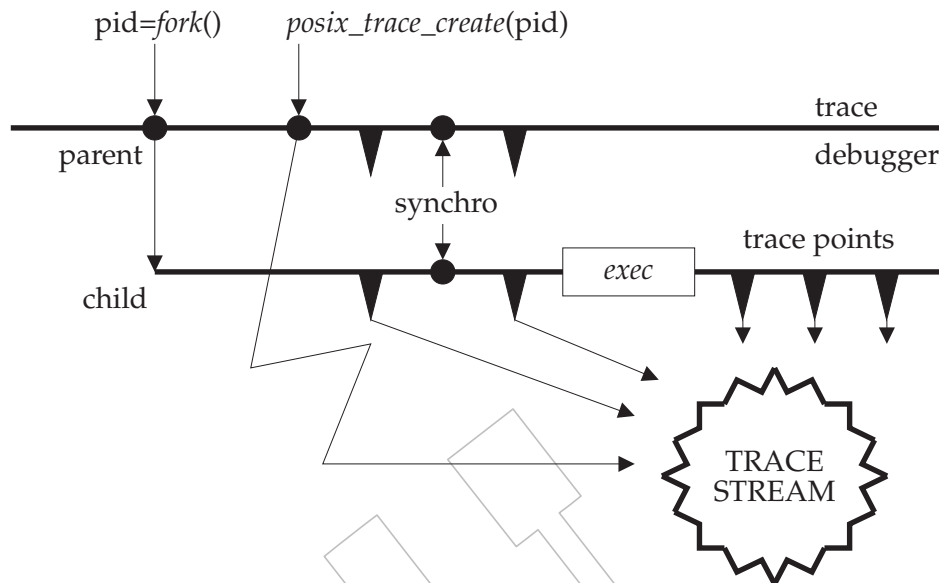
```

8369 /* Caution. Error checks omitted */
8370 {
8371     trace_id_t trid = existing_trace;
8372     trace_event_id_t event_id;
8373     char event_name[TRACE_EVENT_NAME_MAX];
8374     int return_value;
8375     - - - - -
8376     /*
8377      * In a loop print all existing trace event names out
8378      * for the trid trace stream
8379      */
8380     while (1)
8381     {
8382         posix_trace_eventtypelist_getnext_id(trid, &event_id
8383             &return_value);
8384         if (return_value != NULL) break;
8385         /*
8386          * Get the name of the trace event associated
8387          * with trid trace ID
8388          */
8389         posix_trace_eventid_get_name(trid, event_id, event_name);
8390         /* Print the name out */
8391         printf("%s\n", event_name);
8392     }
8393 }

```


8394

B.2.11.4 Rationale on Trace for Debugging



8395

Figure B-5 Trace Another Process

8396

8397

8398

8399

8400

8401

8402

Among the different possibilities offered by the trace interface defined in IEEE Std 1003.1-200x, the debugging of an application is the most interesting one. Typical operations in the controlling debugger process are to filter trace event types, to get trace events from the trace stream, to stop the trace stream when the debugged process is executing uninteresting code, to start the trace stream when some interesting point is reached, and so on. The interface defined in IEEE Std 1003.1-200x should define all the necessary base functions to allow this dynamic debug handling.

8403

8404

8405

8406

Figure B-5 shows an example in which the trace stream is created after the call to the `fork()` function. If the user does not want to lose trace events, some synchronization mechanism (represented in the figure) may be needed before calling the `exec` function, to give the parent a chance to create the trace stream before the child begins the execution of its trace points.

8407

B.2.11.5 Rationale on Trace Event Type Name Space

8408

8409

8410

8411

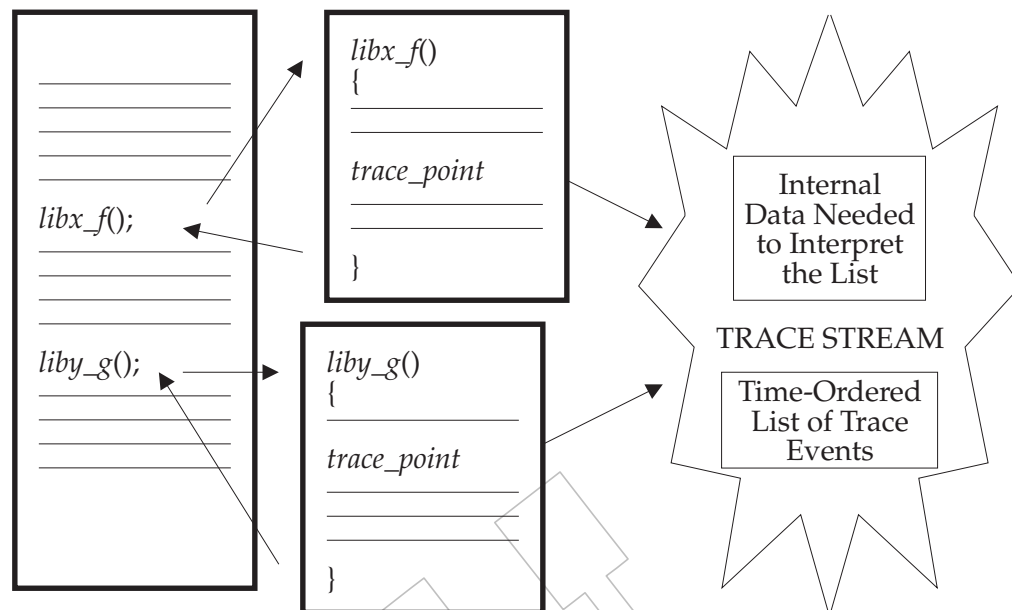
8412

8413

8414

8415

At first, the working group was in favor of the representation of a trace event type by an integer (*event_name*). It seems that existing practice shows the weakness of such a representation. The collision of trace event types is the main problem that cannot be simply resolved using this sort of representation. Suppose, for example, that a third party designs an instrumented library. The user does not have the source of this library and wants to trace his application which uses in some part the third-party library. There is no means for him to know what are the trace event types used in the instrumented library so he has some chance of duplicating some of them and thus to obtain a contaminated tracing of his application.



8416 **Figure B-6** Trace Name Space Overview: With Third-Party Library

8417 There are requirements to allow program images containing pieces from various vendors to be
 8418 traced without also requiring those of any other vendors to coordinate their uses of the trace
 8419 facility, and especially the naming of their various trace event types and trace point IDs. The
 8420 chosen solution is to provide a very large name space, large enough so that the individual
 8421 vendors can give their trace types and tracepoint IDs sufficiently long and descriptive names
 8422 making the occurrence of collisions quite unlikely. The probability of collision is thus made
 8423 sufficiently low so that the problem may, as a practical matter, be ignored. By requirement, the
 8424 consequence of collisions will be a slight ambiguity in the trace streams; tracing will continue in
 8425 spite of collisions and ambiguities. “The show must go on”. The *posix_prog_address* member of
 8426 the **posix_trace_event_info** structure is used to allow trace streams to be unambiguously
 8427 interpreted, despite the fact that trace event types and trace event names need not be unique.

8428 The *posix_trace_eventid_open()* function is required to allow the instrumented third-party library
 8429 to get a valid trace event type identifier for its trace event names. This operation is, somehow,
 8430 an allocation, and the group was aware of proposing some deallocation mechanism which the
 8431 instrumented application could use to recover the resources used by a trace event type identifier.
 8432 This would have given the instrumented application the benefit of being capable of reusing a
 8433 possible minimum set of trace event type identifiers, but also the inconvenience to have,
 8434 possibly in the same trace stream, one trace event type identifier identifying two different trace
 8435 event types. After some discussions the group decided to not define such a function which
 8436 would make this API thicker for little benefit, the user having always the possibility of adding
 8437 identification information in the *data* member of the trace event structure.

8438 The set of the trace event type identifiers the controlling process wants to filter out is initialized
 8439 in the trace mechanism using the function *posix_trace_set_filter()*, setting the arguments
 8440 according to the definitions explained in *posix_trace_set_filter()*. This operation can be done
 8441 statically (when the trace is in the STOPPED state) or dynamically (when the trace is in the
 8442 STARTED state). The preparation of the filter is normally done using the function defined in
 8443 *posix_trace_eventtypelist_getnext_id()* and eventually the function

8444 *posix_trace_eventtypelist_rewind()* in order to know (before the recording) the list of the potential
 8445 set of trace event types that can be recorded. In the case of an active trace stream, this list may
 8446 not be exhaustive. Actually, the target process may not have yet called the function
 8447 *posix_trace_eventid_open()*. But it is a common practice, for a controlling process, to prepare the
 8448 filtering of a future trace stream before its start. Therefore the user must have a way to get the
 8449 trace event type identifier corresponding to a well-known trace event name before its future
 8450 association by the pre-cited function. This is done by calling the *posix_trace_trid_eventid_open()*
 8451 function, given the trace stream identifier and the trace name, and described hereafter. Because
 8452 this trace event type identifier is associated with a trace stream identifier, where a unique
 8453 process has initialized two or more traces, the implementation is expected to return the same
 8454 trace event type identifier for successive calls to *posix_trace_trid_eventid_open()* with different
 8455 trace stream identifiers. The *posix_trace_eventid_get_name()* function is used by the controller
 8456 process to identify, by the name, the trace event type returned by a call to the
 8457 *posix_trace_eventtypelist_getnext_id()* function.

8458 Afterwards, the set of trace event types is constructed using the functions defined in
 8459 *posix_trace_eventset_empty()*, *posix_trace_eventset_fill()*, *posix_trace_eventset_add()*, and
 8460 *posix_trace_eventset_del()*.

8461 A set of functions is provided devoted to the manipulation of the trace event type identifier and
 8462 names for an active trace stream. All these functions require the trace stream identifier argument
 8463 as the first parameter. The opacity of the trace event type identifier implies that the user cannot
 8464 associate directly its well-known trace event name with the system-associated trace event type
 8465 identifier.

8466 The *posix_trace_trid_eventid_open()* function allows the application to get the system trace event
 8467 type identifier back from the system, given its well-known trace event name. This function is
 8468 useful only when a controlling process needs to specify specific events to be filtered.

8469 The *posix_trace_eventid_get_name()* function allows the application to obtain a trace event name
 8470 given its trace event type identifier. One possible use of this function is to identify the type of a
 8471 trace event retrieved from the trace stream, and print it. The easiest way to implement this
 8472 requirement, is to use a single trace event type map for all the processes whose maps are
 8473 required to be identical. A more difficult way is to attempt to keep multiple maps identical at
 8474 every call to *posix_trace_eventid_open()* and *posix_trace_trid_eventid_open()*.

8475 B.2.11.6 Rationale on Trace Events Type Filtering

8476 The most basic rationale for runtime and pre-registration filtering (selection/rejection) of trace
 8477 event types is to prevent choking of the trace collection facility, and/or overloading of the
 8478 computer system. Any worthwhile trace facility can bring even the largest computer to its
 8479 knees. Otherwise, everything would be recorded and filtered after the fact; it would be much
 8480 simpler, but impractical.

8481 To achieve debugging, measurement, or whatever the purpose of tracing, the filtering of trace
 8482 event types is an important part of trace analysis. Due to the fact that the trace events are put
 8483 into a trace stream and probably logged afterwards into a file, different levels of filtering—that
 8484 is, rejection of trace event types—are possible.

8485 **Filtering of Trace Event Types Before Tracing**

8486 This function, represented by the `posix_trace_set_filter()` function in IEEE Std 1003.1-200x (see
 8487 `posix_trace_set_filter()`), selects, before or during tracing, the set of trace event types to be filtered
 8488 out. It should be possible also (as OSF suggested in their ETAP trace specifications) to select the
 8489 kernel trace event types to be traced in a system-wide fashion. These two functionalities are
 8490 called the pre-filtering of trace event types.

8491 The restriction on the actual type used for the `trace_event_set_t` type is intended to guarantee
 8492 that these objects can always be assigned, have their address taken, and be passed by value as
 8493 parameters. It is not intended that this type be a structure including pointers to other data
 8494 structures, as that could impact the portability of applications performing such operations. A
 8495 reasonable implementation could be a structure containing an array of integer types.

8496 **Filtering of Trace Event Types at Runtime**

8497 It is possible to build this functionality using the `posix_trace_set_filter()` function. A privileged
 8498 process or a privileged thread can get trace events from the trace stream of another process or
 8499 thread, and thus specify the type of trace events to record into a file, using implementation-
 8500 defined methods and interfaces. This functionality, called inline filtering of trace event types, is
 8501 used for runtime analysis of trace streams.

8502 **Post-Mortem Filtering of Trace Event Types**

8503 The word “post-mortem” is used here to indicate that some unanticipated situation occurs
 8504 during execution that does not permit a pre or inline filtering of trace events and that it is
 8505 necessary to record all trace event types to have a chance to discover the problem afterwards.
 8506 When the program stops, all the trace events recorded previously can be analyzed in order to
 8507 find the solution. This functionality could be named the post-filtering of trace event types.

8508 **Discussions about Trace Event Type-Filtering**

8509 After long discussions with the parties involved in the process of defining the trace interface, it
 8510 seems that the sensitivity to the filtering problem is different, but everybody agrees that the level
 8511 of the overhead introduced during the tracing operation depends on the filtering method
 8512 elected. If the time that it takes the trace event to be recorded can be neglected, the overhead
 8513 introduced by the filtering process can be classified as follows:

8514	Pre-filtering	System and process/thread-level overhead
8515	Inline-filtering	Process/thread-level overhead
8516	Post-filtering	No overhead; done offline

8517 The pre-filtering could be named “critical realtime” filtering in the sense that the filtering of
 8518 trace event type is manageable at the user level so the user can lower to a minimum the filtering
 8519 overhead at some user selected level of priority for the inline filtering, or delay the filtering to
 8520 after execution for the post-filtering. The counterpart of this solution is that the size of the trace
 8521 stream must be sufficient to record all the trace events. The advantage of the pre-filtering is that
 8522 the utilization of the trace stream is optimized.

8523 Only pre-filtering is defined by IEEE Std 1003.1-200x. However, great care must be taken in
 8524 specifying pre-filtering, so that it does not impose unacceptable overhead. Moreover, it is
 8525 necessary to isolate all the functionality relative to the pre-filtering.

8526 The result of this rationale is to define a new option, the Trace Event Filter option, not
 8527 necessarily implemented in small realtime systems, where system overhead is minimized to the
 8528 extent possible.

8529 B.2.11.7 Tracing, *pthread* API

8530 The objective to be able to control tracing for individual threads may be in conflict with the
 8531 efficiency expected in threads with a *contentionscope* attribute of `PTHREAD_SCOPE_PROCESS`.
 8532 For these threads, context switches from one thread that has tracing enabled to another thread
 8533 that has tracing disabled may require a kernel call to inform the kernel whether it has to trace
 8534 system events executed by that thread or not. For this reason, it was proposed that the ability to
 8535 enable or disable tracing for `PTHREAD_SCOPE_PROCESS` threads be made optional, through
 8536 the introduction of a Trace Scope Process option. A trace implementation which did not
 8537 implement the Trace Scope Process option would not honor the tracing-state attribute of a thread
 8538 with `PTHREAD_SCOPE_PROCESS`; it would, however, honor the tracing-state attribute of a
 8539 thread with `PTHREAD_SCOPE_SYSTEM`. This proposal was rejected as:

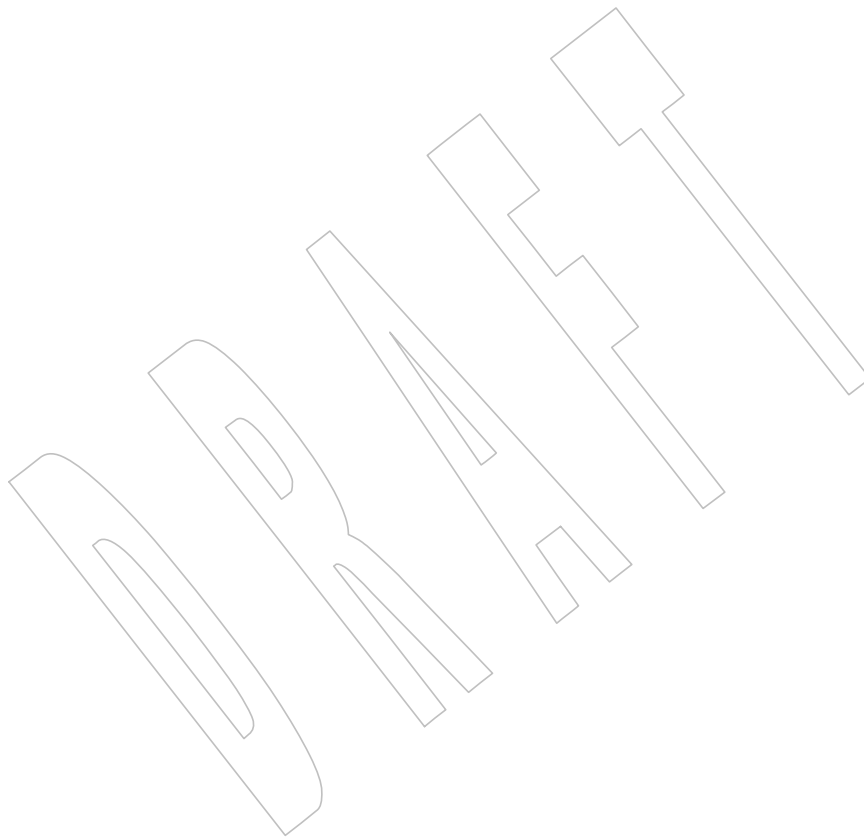
- 8540 1. Removing desired functionality (per-thread trace control)
- 8541 2. Introducing counter-intuitive behavior for the tracing-state attribute
- 8542 3. Mixing logically orthogonal ideas (thread scheduling and thread tracing)
- 8543 [Objective 4]

8544 Finally, to solve this complex issue, this API does not provide `pthread_gettracingstate()`,
 8545 `pthread_settracingstate()`, `pthread_attr_gettracingstate()`, and `pthread_attr_settracingstate()`
 8546 interfaces. These interfaces force the thread implementation to add to the weight of the thread
 8547 and cause a revision of the threads libraries, just to support tracing. Worse yet,
 8548 `posix_trace_event()` must always test this per-thread variable even in the common case where it is
 8549 not used at all. Per-thread tracing is easy to implement using existing interfaces where
 8550 necessary; see the following example.

8551 **Example**

```

8552 /* Caution. Error checks omitted */
8553 static pthread_key_t my_key;
8554 static trace_event_id_t my_event_id;
8555 static pthread_once_t my_once = PTHREAD_ONCE_INIT;
8556
8557 void my_init(void)
8558 {
8559     (void) pthread_key_create(&my_key, NULL);
8560     (void) posix_trace_eventid_open("my", &my_event_id);
8561 }
8562
8563 int get_trace_flag(void)
8564 {
8565     pthread_once(&my_once, my_init);
8566     return (pthread_getspecific(my_key) != NULL);
8567 }
8568
8569 void set_trace_flag(int f)
8570 {
8571     pthread_once(&my_once, my_init);
8572     pthread_setspecific(my_key, f? &my_event_id: NULL);
8573 }
8574
8575 fn()
8576 {
8577     if (get_trace_flag())
8578         posix_trace_event(my_event_id, ...)
8579 }
  
```



8621 **B.2.11.10 Rationale on Different Overrun Conditions**

8622 The analysis of the dynamic behavior of the trace mechanism shows that different overrun
 8623 conditions may occur. The API must provide a means to manage such conditions in a portable
 8624 way.

8625 **Overrun in Trace Streams Initialized with POSIX_TRACE_LOOP Policy**

8626 In this case, the user of the trace mechanism is interested in using the trace stream with
 8627 POSIX_TRACE_LOOP policy to record trace events continuously, but ideally without losing any
 8628 trace events. The online analyzer process must get the trace events at a mean speed equivalent to
 8629 the recording speed. Should the trace stream become full, a trace stream overrun occurs. This
 8630 condition is detected by getting the status of the active trace stream (function
 8631 `posix_trace_get_status()`) and looking at the member `posix_stream_overrun_status` of the read
 8632 `posix_stream_status` structure. In addition, two predefined trace event types are defined:

- 8633 1. The beginning of a trace overflow, to locate the beginning of an overflow when reading a
 8634 trace stream
- 8635 2. The end of a trace overflow, to locate the end of an overflow, when reading a trace stream

8636 As a timestamp is associated with these predefined trace events, it is possible to know the
 8637 duration of the overflow.

8638 **Overrun in Dumping Trace Streams into Trace Logs**

8639 The user lets the trace mechanism dump the trace stream initialized with
 8640 POSIX_TRACE_FLUSH policy automatically into a trace log. If the dump operation is slower
 8641 than the recording of trace events, the trace stream can overrun. This condition is detected by
 8642 getting the status of the active trace stream (the `posix_trace_get_status()` function) and looking at
 8643 the member `posix_stream_overrun_status` of the read `posix_stream_status` structure. This overrun
 8644 indicates that the trace mechanism is not able to operate in this mode at this speed. It is the
 8645 responsibility of the user to modify one of the trace parameters (the stream size or the trace
 8646 event type filter, for instance) to avoid such overrun conditions, if overruns are to be prevented.
 8647 The same already predefined trace event types (see [Overrun in Trace Streams Initialized with
 8648 POSIX_TRACE_LOOP Policy](#) (on page 202)) are used to detect and to know the duration of an
 8649 overflow.

8650 **Reading an Active Trace Stream**

8651 Although this trace API allows one to read an active trace stream with log while it is tracing, this
 8652 feature can lead to false overflow origin interpretation: the trace log or the reader of the trace
 8653 stream. Reading from an active trace stream with log is thus non-portable, and has been left
 8654 unspecified.

8655 **B.2.12 Data Types**8656 **B.2.12.1 Defined Types**

8657 The requirement that additional types defined in this section end in “_t” was prompted by the
 8658 problem of name space pollution. It is difficult to define a type (where that type is not one
 8659 defined by IEEE Std 1003.1-200x) in one header file and use it in another without adding
 8660 symbols to the name space of the program. To allow implementors to provide their own types,
 8661 all conforming applications are required to avoid symbols ending in “_t”, which permits the
 8662 implementor to provide additional types. Because a major use of types is in the definition of
 8663 structure members, which can (and in many cases must) be added to the structures defined in

8664 IEEE Std 1003.1-200x, the need for additional types is compelling.

8665 The types, such as **ushort** and **ulong**, which are in common usage, are not defined in
8666 IEEE Std 1003.1-200x (although **ushort_t** would be permitted as an extension). They can be
8667 added to **<sys/types.h>** using a feature test macro (see Section B.2.2.1 (on page 84)). A
8668 suggested symbol for these is **_SYSIII**. Similarly, the types like **u_short** would probably be best
8669 controlled by **_BSD**.

8670 Some of these symbols may appear in other headers; see Section B.2.2.2 (on page 85).

8671 **dev_t** This type may be made large enough to accommodate host-locality considerations
8672 of networked systems.

8673 This type must be arithmetic. Earlier proposals allowed this to be non-arithmetic
8674 (such as a structure) and provided a *samefile()* function for comparison.

8675 **gid_t** Some implementations had separated **gid_t** from **uid_t** before POSIX.1 was
8676 completed. It would be difficult for them to coalesce them when it was
8677 unnecessary. Additionally, it is quite possible that user IDs might be different from
8678 group IDs because the user ID might wish to span a heterogeneous network,
8679 where the group ID might not.

8680 For current implementations, the cost of having a separate **gid_t** will be only
8681 lexical.

8682 **mode_t** This type was chosen so that implementations could choose the appropriate
8683 integer type, and for compatibility with the ISO C standard. 4.3 BSD uses
8684 **unsigned short** and the SVID uses **ushort**, which is the same. Historically, only the
8685 low-order sixteen bits are significant.

8686 **nlink_t** This type was introduced in place of **short** for *st_nlink* (see the **<sys/stat.h>** header)
8687 in response to an objection that **short** was too small.

8688 **off_t** This type is used only in *lseek()*, *fcntl()*, and **<sys/stat.h>**. Many implementations
8689 would have difficulties if it were defined as anything other than **long**. Requiring
8690 an integer type limits the capabilities of *lseek()* to four gigabytes. The ISO C
8691 standard supplies routines that use larger types; see *fgetpos()* and *fsetpos()*. XSI-
8692 conformant systems provide the *fseeko()* and *ftello()* functions that use larger types.

8693 **pid_t** The inclusion of this symbol was controversial because it is tied to the issue of the
8694 representation of a process ID as a number. From the point of view of a
8695 conforming application, process IDs should be “magic cookies”¹ that are produced
8696 by calls such as *fork()*, used by calls such as *waitpid()* or *kill()*, and not otherwise
8697 analyzed (except that the sign is used as a flag for certain operations).

8698 The concept of a {PID_MAX} value interacted with this in early proposals. Treating
8699 process IDs as an opaque type both removes the requirement for {PID_MAX} and
8700 allows systems to be more flexible in providing process IDs that span a large range
8701 of values, or a small one.

8702 Since the values in **uid_t**, **gid_t**, and **pid_t** will be numbers generally, and
8703 potentially both large in magnitude and sparse, applications that are based on
8704 arrays of objects of this type are unlikely to be fully portable in any case. Solutions
8705 that treat them as magic cookies will be portable.

8706 1. An historical term meaning: “An opaque object, or token, of determinate size, whose significance is known only to the entity which
8707 created it. An entity receiving such a token from the generating entity may only make such use of the ‘cookie’ as is defined and permitted
8708 by the supplying entity.”

8709 {CHILD_MAX} precludes the possibility of a “toy implementation”, where there
8710 would only be one process.

8711 **ssize_t** This is intended to be a signed analog of **size_t**. The wording is such that an
8712 implementation may either choose to use a longer type or simply to use the signed
8713 version of the type that underlies **size_t**. All functions that return **ssize_t** (*read()*
8714 and *write()*) describe as “implementation-defined” the result of an input exceeding
8715 {SSIZE_MAX}. It is recognized that some implementations might have **ints** that
8716 are smaller than **size_t**. A conforming application would be constrained not to
8717 perform I/O in pieces larger than {SSIZE_MAX}, but a conforming application
8718 using extensions would be able to use the full range if the implementation
8719 provided an extended range, while still having a single type-compatible interface.

8720 The symbols **size_t** and **ssize_t** are also required in **<unistd.h>** to minimize the
8721 changes needed for calls to *read()* and *write()*. Implementors are reminded that it
8722 must be possible to include both **<sys/types.h>** and **<unistd.h>** in the same
8723 program (in either order) without error.

8724 **uid_t** Before the addition of this type, the data types used to represent these values
8725 varied throughout early proposals. The **<sys/stat.h>** header defined these values
8726 as type **short**, the **<passwd.h>** file (now **<pwd.h>** and **<grp.h>**) used an **int**, and
8727 *getuid()* returned an **int**. In response to a strong objection to the inconsistent
8728 definitions, all the types were switched to **uid_t**.

8729 In practice, those historical implementations that use varying types of this sort can
8730 typedef **uid_t** to **short** with no serious consequences.

8731 The problem associated with this change concerns object compatibility after
8732 structure size changes. Since most implementations will define **uid_t** as a **short**, the
8733 only substantive change will be a reduction in the size of the **passwd** structure.
8734 Consequently, implementations with an overriding concern for object
8735 compatibility can pad the structure back to its current size. For that reason, this
8736 problem was not considered critical enough to warrant the addition of a separate
8737 type to POSIX.1.

8738 The types **uid_t** and **gid_t** are magic cookies. There is no {UID_MAX} defined by
8739 POSIX.1, and no structure imposed on **uid_t** and **gid_t** other than that they be
8740 positive arithmetic types. (In fact, they could be **unsigned char**.) There is no
8741 maximum or minimum specified for the number of distinct user or group IDs.

8742 B.2.12.2 The char Type

8743 This standard explicitly requires that a **char** type is exactly one byte (8 bits).

8744 B.2.12.3 Pointer Types

8745 This standard explicitly requires implementations to convert pointers to **void *** and back with no
8746 loss of information. This is an extension over the ISO C standard.

8747 B.3 System Interfaces

8748 See the RATIONALE sections on the individual reference pages.

8749 B.3.1 System Interfaces Removed in this Revision

8750 The following section contains a list of the interfaces removed in this revision of the standard,
8751 together with advice for application writers on the alternative interfaces that should be used for
8752 maximum portability.

8753 B.3.1.1 *bcmp*

8754 Applications are recommended to use the *memcmp()* function instead of this function.

8755 For maximum portability, it is recommended to replace the function call to *bcmp()* as follows:

```
8756 #define bcmp(b1,b2,len) memcmp((b1), (b2), (size_t)(len))
```

8757 B.3.1.2 *bcopy*

8758 Applications are recommended to use the *memmove()* function instead of this function.

8759 The following are approximately equivalent (note the order of the arguments):

```
8760 bcopy(s1,s2,n) ~ memmove(s2,s1,n)
```

8761 For maximum portability, it is recommended to replace the function call to *bcopy()* as follows:

```
8762 #define bcopy(b1,b2,len) (void)(memmove((b2), (b1), (len)))
```

8763 B.3.1.3 *bsd_signal*

8764 Applications are recommended to use the *sigaction()* function instead of this function.

8765 The *bsd_signal()* function was supplied as a migration path for the BSD *signal()* function for
8766 simple applications that installed a single-argument signal handler function.

8767 Historically, the *bsd_signal()* function differs from *signal()* in that the SA_RESTART flag is set
8768 and the SA_RESETHAND flag is clear when *bsd_signal()* is used. The state of these flags is not
8769 specified for *signal()*.

8770 B.3.1.4 *bzero*

8771 Applications are recommended to use the *memset()* function instead of this function.

8772 For maximum portability, it is recommended to replace the function call to *bzero()* as follows:

```
8773 #define bzero(b,len) (void)(memset((b), '\0', (len)))
```

8774 B.3.1.5 *ecvt, fcvt, gcvt*

8775 Applications are recommended to use the *sprintf()* function instead of these functions.

8776 The *sprintf()* function is required by ISO C and is thus more portable.

8777 B.3.1.6 *ftime*

8778 Applications are recommended to use the *time()* function to determine the current time.
8779 Realtime applications should use *clock_gettime()* to determine the current time.

8780 B.3.1.7 *getcontext, makecontext, swapcontext*

8781 Due to portability issues with these functions, especially with the manipulation of contexts,
8782 applications are recommended to be rewritten to use POSIX threads.

8783 B.3.1.8 *gethostbyaddr, gethostbyname*

8784 Applications are recommended to use the *getaddrinfo()* and *getnameinfo()* functions instead of
8785 these functions.

8786 The *gethostbyaddr()* and *gethostbyname()* functions may return pointers to static data, which may
8787 be overwritten by subsequent calls to any of these functions. The suggested replacements do not
8788 have this problem and are also IPv6-capable.

8789 B.3.1.9 *getwd*

8790 Applications are recommended to use the *getcwd()* function to determine the current working
8791 directory.

8792 B.3.1.10 *h_errno*

8793 Applications are recommended not to use this error return code. Previously it was set by the
8794 *gethostbyname()* and *gethostbyaddr()* functions.

8795 B.3.1.11 *index*

8796 Applications are recommended to use the *strchr()* function instead of this function.

8797 For maximum portability, it is recommended to replace the function call to *index()* as follows:

```
8798 #define index(a,b) strchr((a),(b))
```

8799 B.3.1.12 *makecontext*

8800 Applications using the *getcontext()*, *makecontext()*, and *swapcontext()* functions should be
8801 rewritten to use POSIX threads.

8802 B.3.1.13 *mktemp*

8803 Applications are recommended to use the *mkstemp()* function instead of this function.

8804 The *mktemp()* function makes an application vulnerable to possible security problems since
8805 between the time a pathname is created and the file opened, it is possible for some other process
8806 to create a file with the same name. The *mkstemp()* function does not have this vulnerability.

8807 B.3.1.14 *pthread_attr_getstackaddr, pthread_attr_setstackaddr*

8808 Applications are recommended to use the *pthread_attr_setstack()* and *pthread_attr_getstack()*
8809 functions instead of these functions.

8810 There are a number of ambiguities in the specification of the *stackaddr* attribute that makes
8811 portable use of these interfaces impossible.

8812 B.3.1.15 *rindex*

8813 Applications are recommended to use the *strrchr()* function instead of this function.

8814 For maximum portability, it is recommended to replace the function call to *rindex()* as follows:

```
8815 #define rindex(a,b) strrchr((a),(b))
```

8816 **B.3.1.16** *scalb*

8817 Applications are recommended to use either *scalbln()*, *scalblnf()*, or *scalblnl()* instead of these
8818 functions.

8819 The behavior for the *scalb()* function was only defined when the *n* argument is an integer, a
8820 NaN, or Inf. The behavior of other values for the *n* argument was unspecified.

8821 **B.3.1.17** *ualarm*

8822 Applications are recommended to use *timer_create()*, *timer_delete()*, *timer_getoverrun()*,
8823 *timer_gettime()*, or *timer_settime()* instead of this function.

8824 **B.3.1.18** *usleep*

8825 Applications are recommended to use the *nanosleep()* function instead of this function.

8826 **B.3.1.19** *vfork*

8827 Applications are recommended to use the *fork()* function instead of this function.

8828 The *vfork()* function was previously under-specified.

8829 **B.3.1.20** *wcs wcs*

8830 Applications are recommended to use the *wcsstr()* function instead of this function.

8831 The *wcsstr()* function is technically equivalent and is portable across all ISO C implementations.

8832 **B.3.2** **Examples for Spawn**

8833 The following long examples are provided in the Rationale (Informative) volume of
8834 IEEE Std 1003.1-200x as a supplement to the reference page for *posix_spawn()*.

8835 **Example Library Implementation of Spawn**

8836 The *posix_spawn()* or *posix_spawnnp()* functions provide the following:

- 8837 • Simply start a process executing a process image. This is the simplest application for
8838 process creation, and it may cover most executions of *fork()*.
- 8839 • Support I/O redirection, including pipes.
- 8840 • Run the child under a user and group ID in the domain of the parent.
- 8841 • Run the child at any priority in the domain of the parent.

8842 The *posix_spawn()* or *posix_spawnnp()* functions do not cover every possible use of the *fork()*
8843 function, but they do span the common applications: typical use by a shell and a login utility.

8844 The price for an application is that before it calls *posix_spawn()* or *posix_spawnnp()*, the parent
8845 must adjust to a state that *posix_spawn()* or *posix_spawnnp()* can map to the desired state for the
8846 child. Environment changes require the parent to save some of its state and restore it afterwards.
8847 The effective behavior of a successful invocation of *posix_spawn()* is as if the operation were
8848 implemented with POSIX operations as follows:

```
8849 #include <sys/types.h>
8850 #include <stdlib.h>
8851 #include <stdio.h>
8852 #include <unistd.h>
8853 #include <sched.h>
8854 #include <fcntl.h>
8855 #include <signal.h>
8856 #include <errno.h>
```



```

8857     #include <string.h>
8858     #include <signal.h>
8859     /* #include <spawn.h> */
8860     /*****
8861     /* Things that could be defined in spawn.h */
8862     /*****
8863     typedef struct
8864     {
8865         short posix_attr_flags;
8866         #define POSIX_SPAWN_SETPGROUP        0x1
8867         #define POSIX_SPAWN_SETSIGMASK      0x2
8868         #define POSIX_SPAWN_SETSIGDEF      0x4
8869         #define POSIX_SPAWN_SETSCHEDULER   0x8
8870         #define POSIX_SPAWN_SETSCHEDPARAM  0x10
8871         #define POSIX_SPAWN_RESETIDS       0x20
8872         pid_t posix_attr_pgroup;
8873         sigset_t posix_attr_sigmask;
8874         sigset_t posix_attr_sigdefault;
8875         int posix_attr_schedpolicy;
8876         struct sched_param posix_attr_schedparam;
8877     } posix_spawnattr_t;
8878     typedef char *posix_spawn_file_actions_t;
8879     int posix_spawn_file_actions_init(
8880     posix_spawn_file_actions_t *file_actions);
8881     int posix_spawn_file_actions_destroy(
8882     posix_spawn_file_actions_t *file_actions);
8883     int posix_spawn_file_actions_addclose(
8884     posix_spawn_file_actions_t *file_actions, int fildes);
8885     int posix_spawn_file_actions_adddup2(
8886     posix_spawn_file_actions_t *file_actions, int fildes,
8887     int newfildes);
8888     int posix_spawn_file_actions_addopen(
8889     posix_spawn_file_actions_t *file_actions, int fildes,
8890     const char *path, int oflag, mode_t mode);
8891     int posix_spawnattr_init(posix_spawnattr_t *attr);
8892     int posix_spawnattr_destroy(posix_spawnattr_t *attr);
8893     int posix_spawnattr_getflags(const posix_spawnattr_t *attr,
8894     short *lags);
8895     int posix_spawnattr_setflags(posix_spawnattr_t *attr, short flags);
8896     int posix_spawnattr_getpgroup(const posix_spawnattr_t *attr,
8897     pid_t *pgroup);
8898     int posix_spawnattr_setpgroup(posix_spawnattr_t *attr, pid_t pgroup);
8899     int posix_spawnattr_getschedpolicy(const posix_spawnattr_t *attr,
8900     int *schedpolicy);
8901     int posix_spawnattr_setschedpolicy(posix_spawnattr_t *attr,
8902     int schedpolicy);
8903     int posix_spawnattr_getschedparam(const posix_spawnattr_t *attr,
8904     struct sched_param *schedparam);
8905     int posix_spawnattr_setschedparam(posix_spawnattr_t *attr,
8906     const struct sched_param *schedparam);
8907     int posix_spawnattr_getsigmask(const posix_spawnattr_t *attr,

```

```

8908     sigset_t *sigmask);
8909 int posix_spawnattr_setsigmask(posix_spawnattr_t *attr,
8910     const sigset_t *sigmask);
8911 int posix_spawnattr_getdefault(const posix_spawnattr_t *attr,
8912     sigset_t *sigdefault);
8913 int posix_spawnattr_setsigdefault(posix_spawnattr_t *attr,
8914     const sigset_t *sigdefault);
8915 int posix_spawn(pid_t *pid, const char *path,
8916     const posix_spawn_file_actions_t *file_actions,
8917     const posix_spawnattr_t *attrp, char *const argv[],
8918     char *const envp[]);
8919 int posix_spawnnp(pid_t *pid, const char *file,
8920     const posix_spawn_file_actions_t *file_actions,
8921     const posix_spawnattr_t *attrp, char *const argv[],
8922     char *const envp[]);
8923
8924 /* Example posix_spawn() library routine */
8925
8926 int posix_spawn(pid_t *pid,
8927     const char *path,
8928     const posix_spawn_file_actions_t *file_actions,
8929     const posix_spawnattr_t *attrp,
8930     char *const argv[],
8931     char *const envp[])
8932 {
8933     /* Create process */
8934     if ((*pid = fork()) == (pid_t) 0)
8935     {
8936         /* This is the child process */
8937         /* Worry about process group */
8938         if (attrp->posix_attr_flags & POSIX_SPAWN_SETPGROUP)
8939         {
8940             /* Override inherited process group */
8941             if (setpgid(0, attrp->posix_attr_pgroup) != 0)
8942             {
8943                 /* Failed */
8944                 exit(127);
8945             }
8946         }
8947
8948         /* Worry about thread signal mask */
8949         if (attrp->posix_attr_flags & POSIX_SPAWN_SETSIGMASK)
8950         {
8951             /* Set the signal mask (can't fail) */
8952             sigprocmask(SIG_SETMASK, &attrp->posix_attr_sigmask, NULL);
8953         }
8954
8955         /* Worry about resetting effective user and group IDs */
8956         if (attrp->posix_attr_flags & POSIX_SPAWN_RESETPGROUP)
8957         {
8958             /* None of these can fail for this case. */
8959             setuid(getuid());
8960             setgid(getgid());

```

```

8959     }
8960     /* Worry about defaulted signals */
8961     if (attrp->posix_attr_flags & POSIX_SPAWN_SETSIGDEF)
8962     {
8963         struct sigaction deflt;
8964         sigset_t all_signals;
8965
8966         int s;
8967
8968         /* Construct default signal action */
8969         deflt.sa_handler = SIG_DFL;
8970         deflt.sa_flags = 0;
8971
8972         /* Construct the set of all signals */
8973         sigfillset(&all_signals);
8974
8975         /* Loop for all signals */
8976         for (s = 0; sigismember(&all_signals, s); s++)
8977         {
8978             /* Signal to be defaulted? */
8979             if (sigismember(&attrp->posix_attr_sigdefault, s))
8980             {
8981                 /* Yes; default this signal */
8982                 if (sigaction(s, &deflt, NULL) == -1)
8983                 {
8984                     /* Failed */
8985                     exit(127);
8986                 }
8987             }
8988         }
8989     }
8990     /* Worry about the fds if they are to be mapped */
8991     if (file_actions != NULL)
8992     {
8993         /* Loop for all actions in object file_actions */
8994         /* (implementation dives beneath abstraction) */
8995         char *p = *file_actions;
8996         while (*p != '\0')
8997         {
8998             if (strncmp(p, "close(", 6) == 0)
8999             {
9000                 int fd;
9001
9002                 if (sscanf(p + 6, "%d", &fd) != 1)
9003                 {
9004                     exit(127);
9005                 }
9006                 if (close(fd) == -1)
9007                     exit(127);
9008             }
9009             else if (strncmp(p, "dup2(", 5) == 0)
9010             {
9011                 int fd, newfd;

```

```

9007         if (sscanf(p + 5, "%d,%d", &fd, &newfd) != 2)
9008         {
9009             exit(127);
9010         }
9011         if (dup2(fd, newfd) == -1)
9012             exit(127);
9013     }
9014     else if (strncmp(p, "open(", 5) == 0)
9015     {
9016         int fd, oflag;
9017         mode_t mode;
9018         int tempfd;
9019         char path[1000];    /* Should be dynamic */
9020         char *q;
9021
9022         if (sscanf(p + 5, "%d,", &fd) != 1)
9023         {
9024             exit(127);
9025         }
9026         p = strchr(p, ',') + 1;
9027         q = strchr(p, '*');
9028         if (q == NULL)
9029             exit(127);
9030         strncpy(path, p, q - p);
9031         path[q - p] = '\\0';
9032         if (sscanf(q + 1, "%o,%o", &oflag, &mode) != 2)
9033         {
9034             exit(127);
9035         }
9036         if (close(fd) == -1)
9037         {
9038             if (errno != EBADF)
9039                 exit(127);
9040         }
9041         tempfd = open(path, oflag, mode);
9042         if (tempfd == -1)
9043             exit(127);
9044         if (tempfd != fd)
9045         {
9046             if (dup2(tempfd, fd) == -1)
9047             {
9048                 exit(127);
9049             }
9050             if (close(tempfd) == -1)
9051             {
9052                 exit(127);
9053             }
9054         }
9055     }
9056     else
9057     {
9058         exit(127);
9059     }

```

```

9059         p = strchr(p, '(') + 1;
9060     }
9061 }
9062 /* Worry about setting new scheduling policy and parameters */
9063 if (attrp->posix_attr_flags & POSIX_SPAWN_SETSCHEDULER)
9064 {
9065     if (sched_setscheduler(0, attrp->posix_attr_schedpolicy,
9066         &attrp->posix_attr_schedparam) == -1)
9067     {
9068         exit(127);
9069     }
9070 }
9071 /* Worry about setting only new scheduling parameters */
9072 if (attrp->posix_attr_flags & POSIX_SPAWN_SETSCHEDPARAM)
9073 {
9074     if (sched_setparam(0, &attrp->posix_attr_schedparam) == -1)
9075     {
9076         exit(127);
9077     }
9078 }
9079 /* Now execute the program at path */
9080 /* Any fd that still has FD_CLOEXEC set will be closed */
9081 execve(path, argv, envp);
9082 exit(127); /* exec failed */
9083 }
9084 else
9085 {
9086     /* This is the parent (calling) process */
9087     if (*pid == (pid_t) - 1)
9088         return errno;
9089     return 0;
9090 }
9091 }
9092
9093 /* *****
9094 /* Here is a crude but effective implementation of the */
9095 /* file action object operators which store actions as */
9096 /* concatenated token-separated strings. */
9097 /* *****
9098 /* Create object with no actions. */
9099 int posix_spawn_file_actions_init(
9100     posix_spawn_file_actions_t *file_actions)
9101 {
9102     *file_actions = malloc(sizeof(char));
9103     if (*file_actions == NULL)
9104         return ENOMEM;
9105     strcpy(*file_actions, "");
9106     return 0;
9107 }
9108 /* Free object storage and make invalid. */
9109 int posix_spawn_file_actions_destroy(

```

```

9109     posix_spawn_file_actions_t *file_actions)
9110 {
9111     free(*file_actions);
9112     *file_actions = NULL;
9113     return 0;
9114 }
9115 /* Add a new action string to object. */
9116 static int add_to_file_actions(
9117     posix_spawn_file_actions_t *file_actions, char *new_action)
9118 {
9119     *file_actions = realloc
9120     (*file_actions, strlen(*file_actions) + strlen(new_action) + 1);
9121     if (*file_actions == NULL)
9122         return ENOMEM;
9123     strcat(*file_actions, new_action);
9124     return 0;
9125 }
9126 /* Add a close action to object. */
9127 int posix_spawn_file_actions_addclose(
9128     posix_spawn_file_actions_t *file_actions, int fildes)
9129 {
9130     char temp[100];
9131     sprintf(temp, "close(%d)", fildes);
9132     return add_to_file_actions(file_actions, temp);
9133 }
9134 /* Add a dup2 action to object. */
9135 int posix_spawn_file_actions_adddup2(
9136     posix_spawn_file_actions_t *file_actions, int fildes,
9137     int newfildes)
9138 {
9139     char temp[100];
9140     sprintf(temp, "dup2(%d,%d)", fildes, newfildes);
9141     return add_to_file_actions(file_actions, temp);
9142 }
9143 /* Add an open action to object. */
9144 int posix_spawn_file_actions_addopen(
9145     posix_spawn_file_actions_t *file_actions, int fildes,
9146     const char *path, int oflag, mode_t mode)
9147 {
9148     char temp[100];
9149     sprintf(temp, "open(%d,%s*%o,%o)", fildes, path, oflag, mode);
9150     return add_to_file_actions(file_actions, temp);
9151 }
9152 /******
9153 /* Here is a crude but effective implementation of the */
9154 /* spawn attributes object functions which manipulate */
9155 /* the individual attributes. */
9156 /******
9157 /* Initialize object with default values. */

```



```

9158     int posix_spawnattr_init(posix_spawnattr_t *attr)
9159     {
9160         attr->posix_attr_flags = 0;
9161         attr->posix_attr_pgroup = 0;
9162         /* Default value of signal mask is the parent's signal mask; */
9163         /* other values are also allowed */
9164         sigprocmask(0, NULL, &attr->posix_attr_sigmask);
9165         sigemptyset(&attr->posix_attr_sigdefault);
9166         /* Default values of scheduling attr inherited from the parent; */
9167         /* other values are also allowed */
9168         attr->posix_attr_schedpolicy = sched_getscheduler(0);
9169         sched_getparam(0, &attr->posix_attr_schedparam);
9170         return 0;
9171     }
9172
9173     int posix_spawnattr_destroy(posix_spawnattr_t *attr)
9174     {
9175         /* No action needed */
9176         return 0;
9177     }
9178
9179     int posix_spawnattr_getflags(const posix_spawnattr_t *attr,
9180     short *flags)
9181     {
9182         *flags = attr->posix_attr_flags;
9183         return 0;
9184     }
9185
9186     int posix_spawnattr_setflags(posix_spawnattr_t *attr, short flags)
9187     {
9188         attr->posix_attr_flags = flags;
9189         return 0;
9190     }
9191
9192     int posix_spawnattr_getpgroup(const posix_spawnattr_t *attr,
9193     pid_t *pgroup)
9194     {
9195         *pgroup = attr->posix_attr_pgroup;
9196         return 0;
9197     }
9198
9199     int posix_spawnattr_setpgroup(posix_spawnattr_t *attr, pid_t pgroup)
9200     {
9201         attr->posix_attr_pgroup = pgroup;
9202         return 0;
9203     }
9204
9205     int posix_spawnattr_getschedpolicy(const posix_spawnattr_t *attr,
9206     int *schedpolicy)
9207     {
9208         *schedpolicy = attr->posix_attr_schedpolicy;
9209         return 0;
9210     }
9211
9212     int posix_spawnattr_setschedpolicy(posix_spawnattr_t *attr,
9213     int schedpolicy)

```

```
9207     {
9208         attr->posix_attr_schedpolicy = schedpolicy;
9209         return 0;
9210     }
9211 int posix_spawnattr_getschedparam(const posix_spawnattr_t *attr,
9212     struct sched_param *schedparam)
9213     {
9214         *schedparam = attr->posix_attr_schedparam;
9215         return 0;
9216     }
9217 int posix_spawnattr_setschedparam(posix_spawnattr_t *attr,
9218     const struct sched_param *schedparam)
9219     {
9220         attr->posix_attr_schedparam = *schedparam;
9221         return 0;
9222     }
9223 int posix_spawnattr_getsigmask(const posix_spawnattr_t *attr,
9224     sigset_t *sigmask)
9225     {
9226         *sigmask = attr->posix_attr_sigmask;
9227         return 0;
9228     }
9229 int posix_spawnattr_setsigmask(posix_spawnattr_t *attr,
9230     const sigset_t *sigmask)
9231     {
9232         attr->posix_attr_sigmask = *sigmask;
9233         return 0;
9234     }
9235 int posix_spawnattr_getsigdefault(const posix_spawnattr_t *attr,
9236     sigset_t *sigdefault)
9237     {
9238         *sigdefault = attr->posix_attr_sigdefault;
9239         return 0;
9240     }
9241 int posix_spawnattr_setsigdefault(posix_spawnattr_t *attr,
9242     const sigset_t *sigdefault)
9243     {
9244         attr->posix_attr_sigdefault = *sigdefault;
9245         return 0;
9246     }
```

9247 **I/O Redirection with Spawn**

9248 I/O redirection with *posix_spawn()* or *posix_spawnnp()* is accomplished by crafting a *file_actions*
9249 argument to effect the desired redirection. Such a redirection follows the general outline of the
9250 following example:

```
9251 /* To redirect new standard output (fd 1) to a file, */  
9252 /* and redirect new standard input (fd 0) from my fd socket_pair[1], */  
9253 /* and close my fd socket_pair[0] in the new process. */  
9254 posix_spawn_file_actions_t file_actions;  
9255 posix_spawn_file_actions_init(&file_actions);  
9256 posix_spawn_file_actions_addopen(&file_actions, 1, "newout", ...);  
9257 posix_spawn_file_actions_dup2(&file_actions, socket_pair[1], 0);  
9258 posix_spawn_file_actions_close(&file_actions, socket_pair[0]);  
9259 posix_spawn_file_actions_close(&file_actions, socket_pair[1]);  
9260 posix_spawn(..., &file_actions, ...);  
9261 posix_spawn_file_actions_destroy(&file_actions);
```

9262 **Spawning a Process Under a New User ID**

9263 Spawning a process under a new user ID follows the outline shown in the following example:

```
9264 Save = getuid();  
9265 setuid(newid);  
9266 posix_spawn(...);  
9267 setuid(Save);
```

9268

/ *Rationale (Informative)*

9269

Part C:

9270

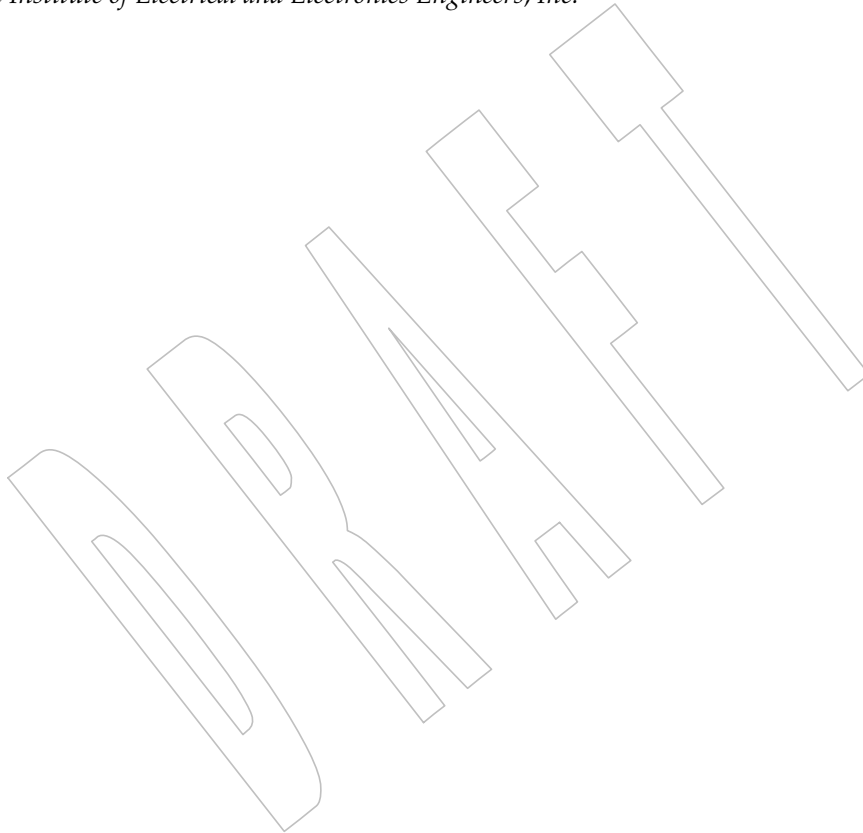
Shell and Utilities

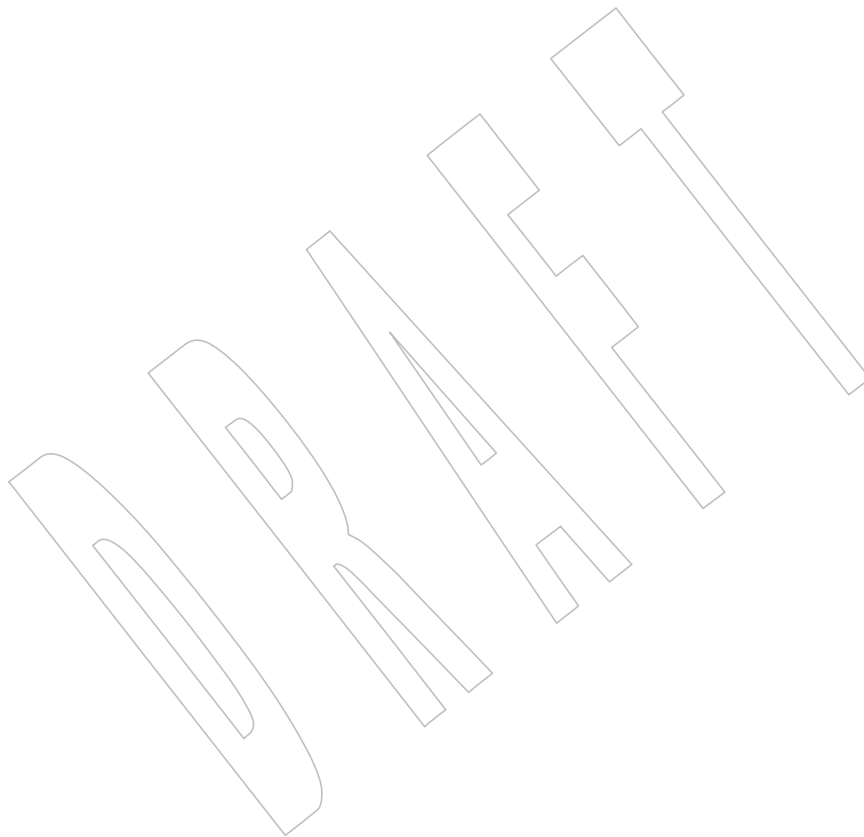
9271

The Open Group

9272

The Institute of Electrical and Electronics Engineers, Inc.





Rationale for Shell and Utilities

9273

9274

9275

C.1 Introduction

9276

C.1.1 Scope

9277

Refer to [Section A.1.1](#) (on page 3).

9278

C.1.2 Conformance

9279

Refer to [Section A.2](#) (on page 9).

9280

C.1.3 Normative References

9281

There is no additional rationale provided for this section.

9282

C.1.4 Change History

9283

The change history is provided as an informative section, to track changes from previous issues of IEEE Std 1003.1-200x.

9284

9285

The following sections describe changes made to the Shell and Utilities volume of IEEE Std 1003.1-200x since Issue 6 of the base document. The CHANGE HISTORY section for each utility describes technical changes made to that utility from Issue 5. Changes between earlier issues of the base document and Issue 5 are not included.

9286

9287

9288

9289

Changes from Issue 6 to Issue 7 (IEEE Std 1003.1-200x)

9290

The following list summarizes the major changes that were made in the Shell and Utilities volume of IEEE Std 1003.1-200x from Issue 6 to Issue 7:

9291

9292

- Austin Group defect reports, IEEE Interpretations against IEEE Std 1003.1, and responses to ISO/IEC defect reports against ISO/IEC 9945 are applied.

9293

9294

- The Open Group corrigenda and resolutions are applied.

9295

- Features, marked legacy or obsolescent in the base document, have been considered for withdrawal in the revision.

9296

9297

- A review of the use of fixed path filenames within the standard has been undertaken; for example, the *at*, *batch*, and *crontab* utilities previously had a requirement for use of the directory **/usr/lib/cron**.

9298

9299

9300

- The options within the standard have been revised.

9301

- The Batch Environment Services and Utilities option is marked obsolescent.

9302

- The UUCP utilities option is added.

9303

- The User Portability Utilities option is revised so that only the *bg*, *ex*, *fc*, *fg*, *jobs*, *more*, *talk*, and *vi* utilities are included, the rest being moved to the Base.

9304

9305 **New Features in Issue 7**

9306 There are no new utilities in Issue 7.

9307 **C.1.5 Terminology**9308 Refer to [Section A.1.4](#) (on page 6).9309 **C.1.6 Definitions**9310 Refer to [Section A.3](#) (on page 13).9311 **C.1.7 Relationship to Other Documents**9312 **C.1.7.1 System Interfaces**

9313 It has been pointed out that the Shell and Utilities volume of IEEE Std 1003.1-200x assumes that
 9314 a great deal of functionality from the System Interfaces volume of IEEE Std 1003.1-200x is
 9315 present, but never states exactly how much (and strictly does not need to since both are
 9316 mandated on a conforming system). This section is an attempt to clarify the assumptions.

9317 **File Read, Write, and Creation**

9318 IEEE Std 1003.1-2001/Cor 2-2004, item XCU/TC2/D6/2 is applied, updating Table 1-1.

9319 **File Removal**

9320 This is intended to be a summary of the *unlink()* and *rmdir()* requirements. Note that it is
 9321 possible using the *unlink()* function for item 4. to occur.

9322 **C.1.7.2 Concepts Derived from the ISO C Standard**

9323 This section was introduced to address the issue that there was insufficient detail presented by
 9324 such utilities as *awk* or *sh* about their procedural control statements and their methods of
 9325 performing arithmetic functions.

9326 The ISO C standard was selected as a model because most historical implementations of the
 9327 standard utilities were written in C. Thus, it was more likely that they would act in the desired
 9328 manner without modification.

9329 Using the ISO C standard is primarily a notational convenience so that the many procedural
 9330 languages in the Shell and Utilities volume of IEEE Std 1003.1-200x would not have to be
 9331 rigorously described in every aspect. Its selection does not require that the standard utilities be
 9332 written in Standard C; they could be written in Common Usage C, Ada, Pascal, assembler
 9333 language, or anything else.

9334 The sizes of the various numeric values refer to C-language data types that are allowed to be
 9335 different sizes by the ISO C standard. Thus, like a C-language application, a shell application
 9336 cannot rely on their exact size. However, it can rely on their minimum sizes expressed in the
 9337 ISO C standard, such as {LONG_MAX} for a **long** type.

9338 The behavior on overflow is undefined for ISO C standard arithmetic. Therefore, the standard
 9339 utilities can use “bignum” representation for integers so that there is no fixed maximum unless
 9340 otherwise stated in the utility description. Similarly, standard utilities can use infinite-precision
 9341 representations for floating-point arithmetic, as long as these representations exceed the ISO C
 9342 standard requirements.

9343 This section addresses only the issue of semantics; it is not intended to specify syntax. For
 9344 example, the ISO C standard requires that 0L be recognized as an integer constant equal to zero,

9345 but utilities such as *awk* and *sh* are not required to recognize 0L (though they are allowed to, as
9346 an extension).

9347 The ISO C standard requires that a C compiler must issue a diagnostic for constants that are too
9348 large to represent. Most standard utilities are not required to issue these diagnostics; for
9349 example, the command:

```
9350 diff -C 2147483648 file1 file2
```

9351 has undefined behavior, and the *diff* utility is not required to issue a diagnostic even if the
9352 number 2 147 483 648 cannot be represented.

9353 C.1.8 Portability

9354 Refer to [Section A.1.5](#) (on page 8).

9355 C.1.8.1 Codes

9356 Refer to [Section A.1.5.1](#) (on page 8).

9357 C.1.9 Utility Limits

9358 This section grew out of an idea that originated with the original POSIX.1, in the tables of system
9359 limits for the *sysconf()* and *pathconf()* functions. The idea being that a conforming application
9360 can be written to use the most restrictive values that a minimal system can provide, but it should
9361 not have to. The values provided represent compromises so that some vendors can use
9362 historically limited versions of UNIX system utilities. They are the highest values that a strictly
9363 conforming application can assume, given no other information.

9364 However, by using the *getconf* utility or the *sysconf()* function, the elegant application can be
9365 tailored to more liberal values on some of the specific instances of specific implementations.

9366 There is no explicitly stated requirement that an implementation provide finite limits for any of
9367 these numeric values; the implementation is free to provide essentially unbounded capabilities
9368 (where it makes sense), stopping only at reasonable points such as {ULONG_MAX} (from the
9369 ISO C standard). Therefore, applications desiring to tailor themselves to the values on a
9370 particular implementation need to be ready for possibly huge values; it may not be a good idea
9371 to allocate blindly a buffer for an input line based on the value of {LINE_MAX}, for instance.
9372 However, unlike the System Interfaces volume of IEEE Std 1003.1-200x, there is no set of limits
9373 that return a special indication meaning “unbounded”. The implementation should always
9374 return an actual number, even if the number is very large.

9375 The statement:

9376 “It is not guaranteed that the application ...”

9377 is an indication that many of these limits are designed to ensure that implementors design their
9378 utilities without arbitrary constraints related to unimaginative programming. There are certainly
9379 conditions under which combinations of options can cause failures that would not render an
9380 implementation non-conforming. For example, {EXPR_NEST_MAX} and {ARG_MAX} could
9381 collide when expressions are large; combinations of {BC_SCALE_MAX} and {BC_DIM_MAX}
9382 could exceed virtual memory.

9383 In the Shell and Utilities volume of IEEE Std 1003.1-200x, the notion of a limit being guaranteed
9384 for the process lifetime, as it is in the System Interfaces volume of IEEE Std 1003.1-200x, is not as
9385 useful to a shell script. The *getconf* utility is probably a process itself, so the guarantee would be
9386 without value. Therefore, the Shell and Utilities volume of IEEE Std 1003.1-200x requires the
9387 guarantee to be for the session lifetime. This will mean that many vendors will either return very
9388 conservative values or possibly implement *getconf* as a built-in.

9389 It may seem confusing to have limits that apply only to a single utility grouped into one global
 9390 section. However, the alternative, which would be to disperse them out into their utility
 9391 description sections, would cause great difficulty when *sysconf()* and *getconf* were described.
 9392 Therefore, the standard developers chose the global approach.

9393 Each language binding could provide symbol names that are slightly different from those shown
 9394 here. For example, the C-Language Binding option adds a leading underscore to the symbols as
 9395 a prefix.

9396 The following comments describe selection criteria for the symbols and their values:

9397 {ARG_MAX}

9398 This is defined by the System Interfaces volume of IEEE Std 1003.1-200x. Unfortunately, it is
 9399 very difficult for a conforming application to deal with this value, as it does not know how
 9400 much of its argument space is being consumed by the environment variables of the user.

9401 {BC_BASE_MAX}

9402 {BC_DIM_MAX}

9403 {BC_SCALE_MAX}

9404 These were originally one value, {BC_SCALE_MAX}, but it was unreasonable to link all
 9405 three concepts into one limit.

9406 {CHILD_MAX}

9407 This is defined by the System Interfaces volume of IEEE Std 1003.1-200x.

9408 {COLL_WEIGHTS_MAX}

9409 The weights assigned to **order** can be considered as “passes” through the collation
 9410 algorithm.

9411 {EXPR_NEST_MAX}

9412 The value for expression nesting was borrowed from the ISO C standard.

9413 {LINE_MAX}

9414 This is a global limit that affects all utilities, unless otherwise noted. The {MAX_CANON}
 9415 value from the System Interfaces volume of IEEE Std 1003.1-200x may further limit input
 9416 lines from terminals. The {LINE_MAX} value was the subject of much debate and is a
 9417 compromise between those who wished to have unlimited lines and those who understood
 9418 that many historical utilities were written with fixed buffers. Frequently, utility writers
 9419 selected the UNIX system constant BUFSIZ to allocate these buffers; therefore, some utilities
 9420 were limited to 512 bytes for I/O lines, while others achieved 4 096 bytes or greater.

9421 It should be noted that {LINE_MAX} applies only to input line length; there is no
 9422 requirement in IEEE Std 1003.1-200x that limits the length of output lines. Utilities such as
 9423 *awk*, *sed*, and *paste* could theoretically construct lines longer than any of the input lines they
 9424 received, depending on the options used or the instructions from the application. They are
 9425 not required to truncate their output to {LINE_MAX}. It is the responsibility of the
 9426 application to deal with this. If the output of one of those utilities is to be piped into another
 9427 of the standard utilities, line length restrictions will have to be considered; the *fold* utility,
 9428 among others, could be used to ensure that only reasonable line lengths reach utilities or
 9429 applications.

9430 {LINK_MAX}

9431 This is defined by the System Interfaces volume of IEEE Std 1003.1-200x.

9432 {MAX_CANON}

9433 {MAX_INPUT}

9434 {NAME_MAX}
 9435 {NGROUPS_MAX}
 9436 {OPEN_MAX}
 9437 {PATH_MAX}
 9438 {PIPE_BUF}

9439 These limits are defined by the System Interfaces volume of IEEE Std 1003.1-200x. Note that
 9440 the byte lengths described by some of these values continue to represent bytes, even if the
 9441 applicable character set uses a multi-byte encoding.

9442 {RE_DUP_MAX}

9443 The value selected is consistent with historical practice. Although the name implies that it
 9444 applies to all REs, only BREs use the interval notation $\{m,n\}$ addressed by this limit.

9445 {POSIX2_SYMLINKS}

9446 The {POSIX2_SYMLINKS} variable indicates that the underlying operating system supports
 9447 the creation of symbolic links in specific directories. Many of the utilities defined in
 9448 IEEE Std 1003.1-200x that deal with symbolic links do not depend on this value. For
 9449 example, a utility that follows symbolic links (or does not, as the case may be) will only be
 9450 affected by a symbolic link if it encounters one. Presumably, a file system that does not
 9451 support symbolic links will not contain any. This variable does affect such utilities as *ln -s*
 9452 and *pax* that attempt to create symbolic links.

9453 There are different limits associated with command lines and input to utilities, depending on the
 9454 method of invocation. In the case of a C program *exec*-ing a utility, {ARG_MAX} is the
 9455 underlying limit. In the case of the shell reading a script and *exec*-ing a utility, {LINE_MAX}
 9456 limits the length of lines the shell is required to process, and {ARG_MAX} will still be a limit. If a
 9457 user is entering a command on a terminal to the shell, requesting that it invoke the utility,
 9458 {MAX_INPUT} may restrict the length of the line that can be given to the shell to a value below
 9459 {LINE_MAX}.

9460 When an option is supported, *getconf* returns a value of 1. For example, when C development is
 9461 supported:

```
9462 if [ "$(getconf POSIX2_C_DEV)" -eq 1 ]; then
9463     echo C supported
9464 fi
```

9465 The *sysconf()* function in the C-Language Binding option would return 1.

9466 The following comments describe selection criteria for the symbols and their values:

```
9467 POSIX2_C_BIND
9468 POSIX2_C_DEV
9469 POSIX2_FORT_DEV
9470 POSIX2_FORT_RUN
9471 POSIX2_SW_DEV
9472 POSIX2_UPE
```

9473 It is possible for some (usually privileged) operations to remove utilities that support these
 9474 options or otherwise to render these options unsupported. The header files, the *sysconf()*
 9475 function, or the *getconf* utility will not necessarily detect such actions, in which case they
 9476 should not be considered as rendering the implementation non-conforming. A test suite
 9477 should not attempt tests such as:

```
9478 rm /usr/bin/c99
9479 getconf POSIX2_C_DEV
```

9480 POSIX2_LOCALEDEF

9481 This symbol was introduced to allow implementations to restrict supported locales to only
9482 those supplied by the implementation.

9483 IEEE Std 1003.1-2001/Cor 1-2002, item XCU/TC1/D6/2 is applied, deleting the entry for
9484 {POSIX2_VERSION} since it is not a utility limit minimum value.

9485 IEEE Std 1003.1-2001/Cor 1-2002, item XCU/TC1/D6/3 is applied, changing the text in Utility
9486 Limits from: “utility (see *getconf*) through the *sysconf*() function defined in the System Interfaces
9487 volume of IEEE Std 1003.1-200x. The literal names shown in Table 1-3 apply only to the *getconf*
9488 utility; the high-level language binding describes the exact form of each name to be used by the
9489 interfaces in that binding.” to: “utility (see *getconf*).”.

9490 C.1.10 Grammar Conventions

9491 There is no additional rationale provided for this section.

9492 C.1.11 Utility Description Defaults

9493 This section is arranged with headings in the same order as all the utility descriptions. It is a
9494 collection of related and unrelated information concerning:

- 9495 1. The default actions of utilities
- 9496 2. The meanings of notations used in IEEE Std 1003.1-200x that are specific to individual
9497 utility sections

9498 Although this material may seem out of place here, it is important that this information appear
9499 before any of the utilities to be described later.

9500 NAME

9501 There is no additional rationale provided for this section.

9502 SYNOPSIS

9503 There is no additional rationale provided for this section.

9504 DESCRIPTION

9505 There is no additional rationale provided for this section.

9506 OPTIONS

9507 Although it has not always been possible, the standard developers tried to avoid repeating
9508 information to reduce the risk that duplicate explanations could each be modified differently.

9509 The need to recognize `--` is required because conforming applications need to shield their
9510 operands from any arbitrary options that the implementation may provide as an extension. For
9511 example, if the standard utility *foo* is listed as taking no options, and the application needed to
9512 give it a pathname with a leading hyphen, it could safely do it as:

9513 `foo -- -myfile`

9514 and avoid any problems with `-m` used as an extension.

9515 **OPERANDS**

9516 The usage of `-` is never shown in the SYNOPSIS. Similarly, the usage of `--` is never shown.

9517 The requirement for processing operands in command-line order is to avoid a “WeirdNIX”
 9518 utility that might choose to sort the input files alphabetically, by size, or by directory order.
 9519 Although this might be acceptable for some utilities, in general the programmer has a right to
 9520 know exactly what order will be chosen.

9521 Some of the standard utilities take multiple *file* operands and act as if they were processing the
 9522 concatenation of those files. For example:

9523 `asa file1 file2`

9524 and:

9525 `cat file1 file2 | asa`

9526 have similar results when questions of file access, errors, and performance are ignored. Other
 9527 utilities such as *grep* or *wc* have completely different results in these two cases. This latter type of
 9528 utility is always identified in its DESCRIPTION or OPERANDS sections, whereas the former is
 9529 not. Although it might be possible to create a general assertion about the former case, the
 9530 following points must be addressed:

- 9531 • Access times for the files might be different in the operand case *versus* the *cat* case.
- 9532 • The utility may have error messages that are cognizant of the input filename, and this
 9533 added value should not be suppressed. (As an example, *awk* sets a variable with the
 9534 filename at each file boundary.)

9535 **STDIN**

9536 There is no additional rationale provided for this section.

9537 **INPUT FILES**

9538 A conforming application cannot assume the following three commands are equivalent:

9539 `tail -n +2 file`
 9540 `(sed -n 1q; cat) < file`
 9541 `cat file | (sed -n 1q; cat)`

9542 The second command is equivalent to the first only when the file is seekable. In the third
 9543 command, if the file offset in the open file description were not unspecified, *sed* would have to
 9544 be implemented so that it read from the pipe 1 byte at a time or it would have to employ some
 9545 method to seek backwards on the pipe. Such functionality is not defined currently in POSIX.1
 9546 and does not exist on all historical systems. Other utilities, such as *head*, *read*, and *sh*, have similar
 9547 properties, so the restriction is described globally in this section.

9548 The definition of “text file” is strictly enforced for input to the standard utilities; very few of
 9549 them list exceptions to the undefined results called for here. (Of course, “undefined” here does
 9550 not mean that historical implementations necessarily have to change to start indicating error
 9551 conditions. Conforming applications cannot rely on implementations succeeding or failing when
 9552 non-text files are used.)

9553 The utilities that allow line continuation are generally those that accept input languages, rather
 9554 than pure data. It would be unusual for an input line of this type to exceed {LINE_MAX} bytes
 9555 and unreasonable to require that the implementation allow unlimited accumulation of multiple
 9556 lines, each of which could reach {LINE_MAX}. Thus, for a conforming application the total of all
 9557 the continued lines in a set cannot exceed {LINE_MAX}.

9558 The format description is intended to be sufficiently rigorous to allow other applications to
 9559 generate these input files. However, since <blank>s can legitimately be included in some of the
 9560 fields described by the standard utilities, particularly in locales other than the POSIX locale, this
 9561 intent is not always realized.

9562 ENVIRONMENT VARIABLES

9563 There is no additional rationale provided for this section.

9564 ASYNCHRONOUS EVENTS

9565 Because there is no language prohibiting it, a utility is permitted to catch a signal, perform some
 9566 additional processing (such as deleting temporary files), restore the default signal action (or
 9567 action inherited from the parent process), and resignal itself.

9568 STDOUT

9569 The format description is intended to be sufficiently rigorous to allow post-processing of output
 9570 by other programs, particularly by an *awk* or *lex* parser.

9571 STDERR

9572 This section does not describe error messages that refer to incorrect operation of the utility.
 9573 Consider a utility that processes program source code as its input. This section is used to
 9574 describe messages produced by a correctly operating utility that encounters an error in the
 9575 program source code on which it is processing. However, a message indicating that the utility
 9576 had insufficient memory in which to operate would not be described.

9577 Some utilities have traditionally produced warning messages without returning a non-zero exit
 9578 status; these are specifically noted in their sections. Other utilities shall not write to standard
 9579 error if they complete successfully, unless the implementation provides some sort of extension to
 9580 increase the verbosity or debugging level.

9581 The format descriptions are intended to be sufficiently rigorous to allow post-processing of
 9582 output by other programs.

9583 OUTPUT FILES

9584 The format description is intended to be sufficiently rigorous to allow post-processing of output
 9585 by other programs, particularly by an *awk* or *lex* parser.

9586 Receipt of the SIGQUIT signal should generally cause termination (unless in some debugging
 9587 mode) that would bypass any attempted recovery actions.

9588 EXTENDED DESCRIPTION

9589 There is no additional rationale provided for this section.

9590 EXIT STATUS

9591 Note the additional discussion of exit values in *Exit Status for Commands* in the *sh* utility. It
 9592 describes requirements for returning exit values greater than 125.

9593 A utility may list zero as a successful return, 1 as a failure for a specific reason, and greater than
 9594 1 as "an error occurred". In this case, unspecified conditions may cause a 2 or 3, or other value,
 9595 to be returned. A strictly conforming application should be written so that it tests for successful
 9596 exit status values (zero in this case), rather than relying upon the single specific error value listed
 9597 in IEEE Std 1003.1-200x. In that way, it will have maximum portability, even on implementations

9598 with extensions.

9599 The standard developers are aware that the general non-enumeration of errors makes it difficult
9600 to write test suites that test the *incorrect* operation of utilities. There are some historical
9601 implementations that have expended effort to provide detailed status messages and a helpful
9602 environment to bypass or explain errors, such as prompting, retrying, or ignoring unimportant
9603 syntax errors; other implementations have not. Since there is no realistic way to mandate system
9604 behavior in cases of undefined application actions or system problems—in a manner acceptable
9605 to all cultures and environments—attention has been limited to the correct operation of utilities
9606 by the conforming application. Furthermore, the conforming application does not need detailed
9607 information concerning errors that it caused through incorrect usage or that it cannot correct.

9608 There is no description of defaults for this section because all of the standard utilities specify
9609 something (or explicitly state “Unspecified”) for exit status.

9610 **CONSEQUENCES OF ERRORS**

9611 Several actions are possible when a utility encounters an error condition, depending on the
9612 severity of the error and the state of the utility. Included in the possible actions of various
9613 utilities are: deletion of temporary or intermediate work files; deletion of incomplete files; and
9614 validity checking of the file system or directory.

9615 The text about recursive traversing is meant to ensure that utilities such as *find* process as many
9616 files in the hierarchy as they can. They should not abandon all of the hierarchy at the first error
9617 and resume with the next command-line operand, but should attempt to keep going.

9618 **APPLICATION USAGE**

9619 This section provides additional caveats, issues, and recommendations to the developer.

9620 **EXAMPLES**

9621 This section provides sample usage.

9622 **RATIONALE**

9623 There is no additional rationale provided for this section.

9624 **FUTURE DIRECTIONS**

9625 FUTURE DIRECTIONS sections act as pointers to related work that may impact the interface in
9626 the future, and often cautions the developer to architect the code to account for a change in this
9627 area. Note that a future directions statement should not be taken as a commitment to adopt a
9628 feature or interface in the future.

9629 **SEE ALSO**

9630 There is no additional rationale provided for this section.

9631 **CHANGE HISTORY**

9632 There is no additional rationale provided for this section.

9633 C.1.12 Considerations for Utilities in Support of Files of Arbitrary Size

9634 This section is intended to clarify the requirements for utilities in support of large files.

9635 The utilities listed in this section are utilities which are used to perform administrative tasks
9636 such as to create, move, copy, remove, change the permissions, or measure the resources of a file.
9637 They are useful both as end-user tools and as utilities invoked by applications during software
9638 installation and operation.

9639 The *chgrp*, *chmod*, *chown*, *ln*, and *rm* utilities probably require use of large file-capable versions of
9640 *stat()*, *lstat()*, *ftw()*, and the **stat** structure.

9641 The *cat*, *cksum*, *cmp*, *cp*, *dd*, *mv*, *sum*, and *touch* utilities probably require use of large file-capable
9642 versions of *creat()*, *open()*, and *fopen()*.

9643 The *cat*, *cksum*, *cmp*, *dd*, *df*, *du*, *ls*, and *sum* utilities may require writing large integer values. For
9644 example:

- 9645 • The *cat* utility might have a `-n` option which counts <newline>s.
- 9646 • The *cksum* and *ls* utilities report file sizes.
- 9647 • The *cmp* utility reports the line number at which the first difference occurs, and also has a
9648 `-l` option which reports file offsets.
- 9649 • The *dd*, *df*, *du*, *ls*, and *sum* utilities report block counts.

9650 The *dd*, *find*, and *test* utilities may need to interpret command arguments that contain 64-bit
9651 values. For *dd*, the arguments include *skip=n*, *seek=n*, and *count=n*. For *find*, the arguments
9652 include `-sizen`. For *test*, the arguments are those associated with algebraic comparisons.

9653 The *df* utility might need to access large file systems with *statvfs()*.

9654 The *ulimit* utility will need to use large file-capable versions of *getrlimit()* and *setrlimit()* and be
9655 able to read and write large integer values.

9656 C.1.13 Built-In Utilities

9657 All of these utilities can be *exec*-ed. There is no requirement that these utilities are actually built
9658 into the shell itself, but many shells need the capability to do so because the Shell and Utilities
9659 volume of IEEE Std 1003.1-200x, Section 2.9.1.1, Command Search and Execution requires that
9660 they be found prior to the *PATH* search. The shell could satisfy its requirements by keeping a list
9661 of the names and directly accessing the file-system versions regardless of *PATH*. Providing all of
9662 the required functionality for those such as *cd* or *read* would be more difficult.

9663 There were originally three justifications for allowing the omission of *exec*-able versions:

- 9664 1. It would require wasting space in the file system, at the expense of very small systems.
9665 However, it has been pointed out that all 16 utilities in the table can be provided with 16
9666 links to a single-line shell script:

```
9667 $0 "$@"
```

- 9668 2. It is not logical to require invocation of utilities such as *cd* because they have no value
9669 outside the shell environment or cannot be useful in a child process. However, counter-
9670 examples always seemed to be available for even the most unusual cases:

```
9671 find . -type d -exec cd {} \; -exec foo {} \;  
9672 (which invokes "foo" on accessible directories)
```

```
9673 ps ... | sed ... | xargs kill
```

```
9674 find . -exec true \; -a ...
```

(where “true” is used for temporary debugging)

3. It is confusing to have a utility such as *kill* that can easily be in the file system in the base standard, but that requires built-in status for the User Portability Utilities option (for the % job control job ID notation). It was decided that it was more appropriate to describe the required functionality (rather than the implementation) to the system implementors and let them decide how to satisfy it.

On the other hand, it was realized that any distinction like this between utilities was not useful to applications, and that the cost to correct it was small. These arguments were ultimately the most effective.

There were varying reasons for including utilities in the table of built-ins:

alias, fc, unalias

The functionality of these utilities is performed more simply within the shell itself and that is the model most historical implementations have used.

bg, fg, jobs

All of the job control-related utilities are eligible for built-in status because that is the model most historical implementations have used.

cd, getopts, newgrp, read, umask, wait

The functionality of these utilities is performed more simply within the context of the current process. An example can be taken from the usage of the *cd* utility. The purpose of the *cd* utility is to change the working directory for subsequent operations. The actions of *cd* affect the process in which *cd* is executed and all subsequent child processes of that process. Based on the POSIX standard process model, changes in the process environment of a child process have no effect on the parent process. If the *cd* utility were executed from a child process, the working directory change would be effective only in the child process. Child processes initiated subsequent to the child process that executed the *cd* utility would not have a changed working directory relative to the parent process.

command

This utility was placed in the table primarily to protect scripts that are concerned about their *PATH* being manipulated. The “secure” shell script example in the *command* utility in the Shell and Utilities volume of IEEE Std 1003.1-200x would not be possible if a *PATH* change retrieved an alien version of *command*. (An alternative would have been to implement *getconf* as a built-in, but the standard developers considered that it carried too many changing configuration strings to require in the shell.)

kill Since *kill* provides optional job control functionality using shell notation (%1, %2, and so on), some implementations would find it extremely difficult to provide this outside the shell.

true, false

These are in the table as a courtesy to programmers who wish to use the “while true” shell construct without protecting *true* from *PATH* searches. (It is acknowledged that “while :” also works, but the idiom with *true* is historically pervasive.)

All utilities, including those in the table, are accessible via the *system()* and *popen()* functions in the System Interfaces volume of IEEE Std 1003.1-200x. There are situations where the return functionality of *system()* and *popen()* is not desirable. Applications that require the exit status of the invoked utility will not be able to use *system()* or *popen()*, since the exit status returned is that of the command language interpreter rather than that of the invoked utility. The alternative for such applications is the use of the *exec* family.

9720 C.2 Shell Command Language

9721 C.2.1 Shell Introduction

9722 The System V shell was selected as the starting point for the Shell and Utilities volume of
9723 IEEE Std 1003.1-200x. The BSD C shell was excluded from consideration for the following
9724 reasons:

- 9725 • Most historically portable shell scripts assume the Version 7 Bourne shell, from which the
9726 System V shell is derived.
- 9727 • The majority of tutorial materials on shell programming assume the System V shell.

9728 The construct "#!" is reserved for implementations wishing to provide that extension. If it were
9729 not reserved, the Shell and Utilities volume of IEEE Std 1003.1-200x would disallow it by forcing
9730 it to be a comment. As it stands, a strictly conforming application must not use "#!" as the first
9731 two characters of the file.

9732 C.2.2 Quoting

9733 There is no additional rationale provided for this section.

9734 C.2.2.1 Escape Character (Backslash)

9735 There is no additional rationale provided for this section.

9736 C.2.2.2 Single-Quotes

9737 A backslash cannot be used to escape a single-quote in a single-quoted string. An embedded
9738 quote can be created by writing, for example: "'a'\''b'", which yields "a'b". (See the Shell
9739 and Utilities volume of IEEE Std 1003.1-200x, Section 2.6.5, Field Splitting for a better
9740 understanding of how portions of words are either split into fields or remain concatenated.) A
9741 single token can be made up of concatenated partial strings containing all three kinds of quoting
9742 or escaping, thus permitting any combination of characters.

9743 C.2.2.3 Double-Quotes

9744 The escaped <newline> used for line continuation is removed entirely from the input and is not
9745 replaced by any white space. Therefore, it cannot serve as a token separator.

9746 In double-quoting, if a backslash is immediately followed by a character that would be
9747 interpreted as having a special meaning, the backslash is deleted and the subsequent character is
9748 taken literally. If a backslash does not precede a character that would have a special meaning, it
9749 is left in place unmodified and the character immediately following it is also left unmodified.
9750 Thus, for example:

9751 "\\$" → \$

9752 "\a" → \a

9753 It would be desirable to include the statement "The characters from an enclosed "\${" to the
9754 matching '}' shall not be affected by the double quotes", similar to the one for "\$()".
9755 However, historical practice in the System V shell prevents this.

9756 The requirement that double-quotes be matched inside "\${...}" within double-quotes and the
9757 rule for finding the matching '}' in the Shell and Utilities volume of IEEE Std 1003.1-200x,
9758 Section 2.6.2, Parameter Expansion eliminate several subtle inconsistencies in expansion for
9759 historical shells in rare cases; for example:

9760 "\${foo-bar }

9761 yields **bar** when **foo** is not defined, and is an invalid substitution when **foo** is defined, in many
 9762 historical shells. The differences in processing the "\${...}" form have led to inconsistencies
 9763 between historical systems. A consequence of this rule is that single-quotes cannot be used to
 9764 quote the '}' within "\${...}"; for example:

```
9765 unset bar
9766 foo="${bar-'}'"
```

9767 is invalid because the "\${...}" substitution contains an unpaired unescaped single-quote. The
 9768 backslash can be used to escape the '}' in this example to achieve the desired result:

```
9769 unset bar
9770 foo="${bar-\}]"
```

9771 The differences in processing the "\${...}" form have led to inconsistencies between the
 9772 historical System V shell, BSD, and KornShells, and the text in the Shell and Utilities volume of
 9773 IEEE Std 1003.1-200x is an attempt to converge them without breaking too many applications.
 9774 The only alternative to this compromise between shells would be to make the behavior
 9775 unspecified whenever the literal characters ' ', '{', '}', and ' "' appear within "\${...}".
 9776 To write a portable script that uses these values, a user would have to assign variables; for
 9777 example:

```
9778 squote=\` dquote=\" lbrace='{` rbrace='}'
9779 ${foo-$squote$rbrace$squote}
```

9780 rather than:

```
9781 ${foo-" ' ' " }
```

9782 Some implementations have allowed the end of the word to terminate the backquoted command
 9783 substitution, such as in:

```
9784 "`echo hello"
```

9785 This usage is undefined; the matching backquote is required by the Shell and Utilities volume of
 9786 IEEE Std 1003.1-200x. The other undefined usage can be illustrated by the example:

```
9787 sh -c '`echo "foo`'
```

9788 The description of the recursive actions involving command substitution can be illustrated with
 9789 an example. Upon recognizing the introduction of command substitution, the shell parses input
 9790 (in a new context), gathering the source for the command substitution until an unbalanced ')' or
 9791 or '`' is located. For example, in the following:

```
9792 echo "$(date; echo "  

  9793     one" )"
```

9794 the double-quote following the *echo* does not terminate the first double-quote; it is part of the
 9795 command substitution script. Similarly, in:

```
9796 echo "$(echo *)"
```

9797 the asterisk is not quoted since it is inside command substitution; however:

```
9798 echo "$(echo "*" )"
```

9799 is quoted (and represents the asterisk character itself).

9800 C.2.3 Token Recognition

9801 The "(" and ")" symbols are control operators in the KornShell, used for an alternative
 9802 syntax of an arithmetic expression command. A conforming application cannot use "(" as a
 9803 single token (with the exception of the "\$(" form for shell arithmetic).

9804 On some implementations, the symbol "(" is a control operator; its use produces unspecified
 9805 results. Applications that wish to have nested subshells, such as:

```
9806 ((echo Hello);(echo World))
```

9807 must separate the "(" characters into two tokens by including white space between them.
 9808 Some systems may treat these as invalid arithmetic expressions instead of subshells.

9809 Certain combinations of characters are invalid in portable scripts, as shown in the grammar.
 9810 Implementations may use these combinations (such as "|&") as valid control operators. Portable
 9811 scripts cannot rely on receiving errors in all cases where this volume of IEEE Std 1003.1-200x
 9812 indicates that a syntax is invalid.

9813 The (3) rule about combining characters to form operators is not meant to preclude systems from
 9814 extending the shell language when characters are combined in otherwise invalid ways.
 9815 Conforming applications cannot use invalid combinations, and test suites should not penalize
 9816 systems that take advantage of this fact. For example, the unquoted combination "|&" is not
 9817 valid in a POSIX script, but has a specific KornShell meaning.

9818 The (10) rule about '#' as the current character is the first in the sequence in which a new token
 9819 is being assembled. The '#' starts a comment only when it is at the beginning of a token. This
 9820 rule is also written to indicate that the search for the end-of-comment does not consider escaped
 9821 <newline> specially, so that a comment cannot be continued to the next line.

9822 C.2.3.1 Alias Substitution

9823 The alias capability was added because it is widely used in historical implementations by
 9824 interactive users.

9825 The definition of "alias name" precludes an alias name containing a slash character. Since the
 9826 text applies to the command words of simple commands, reserved words (in their proper
 9827 places) cannot be confused with aliases.

9828 The placement of alias substitution in token recognition makes it clear that it precedes all of the
 9829 word expansion steps.

9830 An example concerning trailing <blank>s and reserved words follows. If the user types:

```
9831 $ alias foo="/bin/ls "  

  9832 $ alias while="/"
```

9833 The effect of executing:

```
9834 $ while true  

  9835 > do  

  9836 > echo "Hello, World"  

  9837 > done
```

9838 is a never-ending sequence of "Hello, World" strings to the screen. However, if the user
 9839 types:

```
9840 $ foo while
```

9841 the result is an *ls* listing of /. Since the alias substitution for **foo** ends in a <space>, the next word
 9842 is checked for alias substitution. The next word, **while**, has also been aliased, so it is substituted
 9843 as well. Since it is not in the proper position as a command word, it is not recognized as a

9844 reserved word.
 9845 If the user types:
 9846 `$ foo; while`
 9847 `while` retains its normal reserved-word properties.

9848 C.2.4 Reserved Words

9849 All reserved words are recognized syntactically as such in the contexts described. However, note
 9850 that `in` is the only meaningful reserved word after a `case` or `for`; similarly, `in` is not meaningful as
 9851 the first word of a simple command.

9852 Reserved words are recognized only when they are delimited (that is, meet the definition of the
 9853 Base Definitions volume of IEEE Std 1003.1-200x, Section 3.435, Word), whereas operators are
 9854 themselves delimiters. For instance, `'(' and ')'` are control operators, so that no `<space>`
 9855 is needed in *(list)*. However, `'{ ' and '}'` are reserved words in *{list}*; so that in this case the
 9856 leading `<space>` and semicolon are required.

9857 The list of unspecified reserved words is from the KornShell, so conforming applications cannot
 9858 use them in places a reserved word would be recognized. This list contained `time` in early
 9859 proposals, but it was removed when the *time* utility was selected for the Shell and Utilities
 9860 volume of IEEE Std 1003.1-200x.

9861 There was a strong argument for promoting braces to operators (instead of reserved words), so
 9862 they would be syntactically equivalent to subshell operators. Concerns about compatibility
 9863 outweighed the advantages of this approach. Nevertheless, conforming applications should
 9864 consider quoting `'{ ' and '}'` when they represent themselves.

9865 The restriction on ending a name with a colon is to allow future implementations that support
 9866 named labels for flow control; see the RATIONALE for the *break* built-in utility.

9867 It is possible that a future version of the Shell and Utilities volume of IEEE Std 1003.1-200x may
 9868 require that `'{ ' and '}'` be treated individually as control operators, although the token `"{ }"`
 9869 will probably be a special-case exemption from this because of the often-used *find{ }* construct.

9870 C.2.5 Parameters and Variables

9871 C.2.5.1 Positional Parameters

9872 There is no additional rationale provided for this section.

9873 C.2.5.2 Special Parameters

9874 Most historical implementations implement subshells by forking; thus, the special parameter
 9875 `'$ ' does not necessarily represent the process ID of the shell process executing the commands`
 9876 `since the subshell execution environment preserves the value of '$ '.`

9877 If a subshell were to execute a background command, the value of `"$!"` for the parent would
 9878 not change. For example:

```
9879 (
9880 date &
9881 echo $!
9882 )
9883 echo $!
```

9884 would echo two different values for `"$!"`.

9885 The "\$-" special parameter can be used to save and restore *set* options:

```
9886 Save=$(echo $- | sed 's/[ics]//g')
9887 ...
9888 set +aCefnuvx
9889 if [ -n "$Save" ]; then
9890     set -$Save
9891 fi
```

9892 The three options are removed using *sed* in the example because they may appear in the value of
9893 "\$-" (from the *sh* command line), but are not valid options to *set*.

9894 The descriptions of parameters '*' and '@' assume the reader is familiar with the field
9895 splitting discussion in the Shell and Utilities volume of IEEE Std 1003.1-200x, Section 2.6.5, Field
9896 Splitting and understands that portions of the word remain concatenated unless there is some
9897 reason to split them into separate fields.

9898 Some examples of the '*' and '@' properties, including the concatenation aspects:

```
9899 set "abc" "def ghi" "jkl"
9900 echo $*      => "abc" "def" "ghi" "jkl"
9901 echo "$*"    => "abc def ghi jkl"
9902 echo $@      => "abc" "def" "ghi" "jkl"
9903 but:
9904 echo "$@"    => "abc" "def ghi" "jkl"
9905 echo "xx$@yy" => "xxabc" "def ghi" "jkl"yy"
9906 echo "$@$@"  => "abc" "def ghi" "jklabc" "def ghi" "jkl"
```

9907 In the preceding examples, the double-quote characters that appear after the "=>" do not
9908 appear in the output and are used only to illustrate word boundaries.

9909 The following example illustrates the effect of setting *IFS* to a null string:

```
9910 $ IFS= ''
9911 $ set foo bar bam
9912 $ echo "$@"
9913 foo bar bam
9914 $ echo "$*"
9915 foobarbam
9916 $ unset IFS
9917 $ echo "$*"
9918 foo bar bam
```

9919 C.2.5.3 Shell Variables

9920 See the discussion of *IFS* in [Section C.2.6.5](#) and the RATIONALE for the *sh* utility.

9921 The prohibition on *LC_CTYPE* changes affecting lexical processing protects the shell
9922 implementor (and the shell programmer) from the ill effects of changing the definition of
9923 <blank> or the set of alphabetic characters in the current environment. It would probably not be
9924 feasible to write a compiled version of a shell script without this rule. The rule applies only to
9925 the current invocation of the shell and its subshells—invoking a shell script or performing *exec*
9926 *sh* would subject the new shell to the changes in *LC_CTYPE*.

9927 Other common environment variables used by historical shells are not specified by the Shell and
9928 Utilities volume of IEEE Std 1003.1-200x, but they should be reserved for the historical uses.

9929 Tilde expansion for components of *PATH* in an assignment such as:

9930	<code>PATH=~hlj/bin:~dwc/bin:\$PATH</code>	
9931		is a feature of some historical shells and is allowed by the wording of the Shell and Utilities
9932		volume of IEEE Std 1003.1-200x, Section 2.6.1, Tilde Expansion. Note that the tildes are
9933		expanded during the assignment to <i>PATH</i> , not when <i>PATH</i> is accessed during command search.
9934		The following entries represent additional information about variables included in the Shell and
9935		Utilities volume of IEEE Std 1003.1-200x, or rationale for common variables in use by shells that
9936		have been excluded:
9937	—	(Underscore.) While underscore is historical practice, its overloaded usage in
9938		the KornShell is confusing, and it has been omitted from the Shell and Utilities
9939		volume of IEEE Std 1003.1-200x.
9940	<i>ENV</i>	This variable can be used to set aliases and other items local to the invocation
9941		of a shell. The file referred to by <i>ENV</i> differs from <i>\$HOME/.profile</i> in that
9942		<i>.profile</i> is typically executed at session start-up, whereas the <i>ENV</i> file is
9943		executed at the beginning of each shell invocation. The <i>ENV</i> value is
9944		interpreted in a manner similar to a dot script, in that the commands are
9945		executed in the current environment and the file needs to be readable, but not
9946		executable. However, unlike dot scripts, no <i>PATH</i> searching is performed. This
9947		is used as a guard against Trojan Horse security breaches.
9948	<i>ERRNO</i>	This variable was omitted from the Shell and Utilities volume of
9949		IEEE Std 1003.1-200x because the values of error numbers are not defined in
9950		IEEE Std 1003.1-200x in a portable manner.
9951	<i>FCEDIT</i>	Since this variable affects only the <i>fc</i> utility, it has been omitted from this more
9952		global place. The value of <i>FCEDIT</i> does not affect the command-line editing
9953		mode in the shell; see the description of <i>set -o vi</i> in the <i>set</i> built-in utility.
9954	<i>PS1</i>	This variable is used for interactive prompts. Historically, the “superuser”
9955		has had a prompt of ‘#’. Since privileges are not required to be monolithic, it
9956		is difficult to define which privileges should cause the alternate prompt.
9957		However, a sufficiently powerful user should be reminded of that power by
9958		having an alternate prompt.
9959	<i>PS3</i>	This variable is used by the KornShell for the <i>select</i> command. Since the POSIX
9960		shell does not include <i>select</i> , <i>PS3</i> was omitted.
9961	<i>PS4</i>	This variable is used for shell debugging. For example, the following script:
9962		<code>PS4='[\${LINENO}]+ '</code>
9963		<code>set -x</code>
9964		<code>echo Hello</code>
9965		writes the following to standard error:
9966		<code>[3]+ echo Hello</code>
9967	<i>RANDOM</i>	This pseudo-random number generator was not seen as being useful to
9968		interactive users.
9969	<i>SECONDS</i>	Although this variable is sometimes used with <i>PS1</i> to allow the display of the
9970		current time in the prompt of the user, it is not one that would be manipulated
9971		frequently enough by an interactive user to include in the Shell and Utilities
9972		volume of IEEE Std 1003.1-200x.

9973

C.2.6 Word Expansions

9974

9975

9976

Step (2) refers to the “portions of fields generated by step (1)”. For example, if the word being expanded were “\$x+\$y” and *IFS*=+, the word would be split only if “\$x” or “\$y” contained ‘+’; the ‘+’ in the original word was not generated by step (1).

9977

9978

9979

9980

9981

9982

IFS is used for performing field splitting on the results of parameter and command substitution; it is not used for splitting all fields. Previous versions of the shell used it for splitting all fields during field splitting, but this has severe problems because the shell can no longer parse its own script. There are also important security implications caused by this behavior. All useful applications of *IFS* use it for parsing input of the *read* utility and for splitting the results of parameter and command substitution.

9983

The rule concerning expansion to a single field requires that if **foo=abc** and **bar=def**, that:

9984

```
"$foo" "$bar"
```

9985

expands to the single field:

9986

```
abcdef
```

9987

The rule concerning empty fields can be illustrated by:

9988

```
$ unset foo
```

9989

```
$ set $foo bar ' ' xyz "$foo" abc
```

9990

```
$ for i
```

9991

```
> do
```

9992

```
>     echo "-$i-"
```

9993

```
> done
```

9994

```
-bar-
```

9995

```
--
```

9996

```
-xyz-
```

9997

```
--
```

9998

```
-abc-
```

9999

Step (1) indicates that parameter expansion, command substitution, and arithmetic expansion are all processed simultaneously as they are scanned. For example, the following is valid arithmetic:

10000

10001

10002

```
x=1
```

10003

```
echo $(( $(echo 3)+$x ))
```

10004

An early proposal stated that tilde expansion preceded the other steps, but this is not the case in known historical implementations; if it were, and if a referenced home directory contained a ‘\$’ character, expansions would result within the directory name.

10005

10006

10007

C.2.6.1 Tilde Expansion

10008

Tilde expansion generally occurs only at the beginning of words, but an exception based on historical practice has been included:

10009

10010

```
PATH=/posix/bin:~dgk/bin
```

10011

This is eligible for tilde expansion because tilde follows a colon and none of the relevant characters is quoted. Consideration was given to prohibiting this behavior because any of the following are reasonable substitutes:

10012

10013

10014

```
PATH=$(printf %s ~karels/bin : ~bostic/bin)
```

10015

```
for Dir in ~maat/bin ~srb/bin ...
```

10016

```
do
```

10017 PATH=\${PATH:+\$PATH:}\$Dir
10018 done

10019 In the first command, explicit colons are used for each directory. In all cases, the shell performs
10020 tilde expansion on each directory because all are separate words to the shell.

10021 Note that expressions in operands such as:

10022 make -k mumble LIBDIR=~chet/lib

10023 do not qualify as shell variable assignments, and tilde expansion is not performed (unless the
10024 command does so itself, which *make* does not).

10025 Because of the requirement that the word is not quoted, the following are not equivalent; only
10026 the last causes tilde expansion:

10027 \~hlj/ ~h\lj/ ~"hlj"/ ~hlj\ ~hlj/

10028 In an early proposal, tilde expansion occurred following any unquoted equals sign or colon, but
10029 this was removed because of its complexity and to avoid breaking commands such as:

10030 rcp hostname:~marc/.profile .

10031 A suggestion was made that the special sequence "\$~" should be allowed to force tilde
10032 expansion anywhere. Since this is not historical practice, it has been left for future
10033 implementations to evaluate. (The description in the Shell and Utilities volume of
10034 IEEE Std 1003.1-200x, Section 2.2, Quoting requires that a dollar sign be quoted to represent
10035 itself, so the "\$~" combination is already unspecified.)

10036 The results of giving tilde with an unknown login name are undefined because the KornShell
10037 "~+" and "~-" constructs make use of this condition, but in general it is an error to give an
10038 incorrect login name with tilde. The results of having *HOME* unset are unspecified because
10039 some historical shells treat this as an error.

10040 C.2.6.2 Parameter Expansion

10041 The rule for finding the closing '}' in "\${...}" is the one used in the KornShell and is
10042 upwardly-compatible with the Bourne shell, which does not determine the closing '}' until the
10043 word is expanded. The advantage of this is that incomplete expansions, such as:

10044 \${foo

10045 can be determined during tokenization, rather than during expansion.

10046 The string length and substring capabilities were included because of the demonstrated need for
10047 them, based on their usage in other shells, such as C shell and KornShell.

10048 Historical versions of the KornShell have not performed tilde expansion on the word part of
10049 parameter expansion; however, it is more consistent to do so.

10050 C.2.6.3 Command Substitution

10051 The "\$()" form of command substitution solves a problem of inconsistent behavior when using
10052 backquotes. For example:

Command	Output
echo '\\$x'	\\$x
echo `echo '\\$x'`	\$x
echo \$(echo '\\$x')	\\$x

10057 Additionally, the backquoted syntax has historical restrictions on the contents of the embedded
10058 command. While the newer "\$()" form can process any kind of valid embedded script, the

10059 backquoted form cannot handle some valid scripts that include backquotes. For example, these
 10060 otherwise valid embedded scripts do not work in the left column, but do work on the right:

10061	echo `	echo \$(
10062	cat <<\eof	cat <<\eof
10063	a here-doc with `	a here-doc with)
10064	eof	eof
10065	`)
10066	echo `	echo \$(
10067	echo abc # a comment with `	echo abc # a comment with)
10068	`)
10069	echo `	echo \$(
10070	echo ` ` `	echo ` ` `
10071	`)

10072 Because of these inconsistent behaviors, the backquoted variety of command substitution is not
 10073 recommended for new applications that nest command substitutions or attempt to embed
 10074 complex scripts.

10075 The KornShell feature:

10076 If *command* is of the form `<word>`, *word* is expanded to generate a pathname, and the value
 10077 of the command substitution is the contents of this file with any trailing `<newline>`
 10078 deleted.

10079 was omitted from the Shell and Utilities volume of IEEE Std 1003.1-200x because `$(cat word)` is
 10080 an appropriate substitute. However, to prevent breaking numerous scripts relying on this
 10081 feature, it is unspecified to have a script within `"$()"` that has only redirections.

10082 The requirement to separate `"$("` and `'('` when a single subshell is command-substituted is to
 10083 avoid any ambiguities with arithmetic expansion.

10084 IEEE Std 1003.1-2001/Cor 1-2002, item XCU/TC1/D6/4 is applied, changing the text from: "If a
 10085 command substitution occurs inside double-quotes, it shall not be performed on the results of
 10086 the substitution." to: "If a command substitution occurs inside double-quotes, field splitting and
 10087 pathname expansion shall not be performed on the results of the substitution.". The
 10088 replacement text taken from the ISO POSIX-2:1993 standard is clearer about the items that are
 10089 not performed.

10090 SD5-XCU-ERN-84 is applied, clarifying how the search for the matching backquote is satisfied.

10091 C.2.6.4 Arithmetic Expansion

10092 The standard developers agreed that there was a strong desire for some kind of arithmetic
 10093 evaluator to provide functionality similar to *expr*, that relating it to `'$'` makes it work well with
 10094 the standard shell language and provides access to arithmetic evaluation in places where
 10095 accessing a utility would be inconvenient.

10096 The syntax and semantics for arithmetic were revised for the ISO/IEC 9945-2:1993 standard.
 10097 The language represents a simple subset of the previous arithmetic language (which was
 10098 derived from the KornShell `"(())"` construct). The syntax was changed from that of a
 10099 command denoted by `((expression))` to an expansion denoted by `$(expression)`. The new form is
 10100 a dollar expansion `'$'` that evaluates the expression and substitutes the resulting value.
 10101 Objections to the previous style of arithmetic included that it was too complicated, did not fit in
 10102 well with the use of variables in the shell, and its syntax conflicted with subshells. The
 10103 justification for the new syntax is that the shell is traditionally a macro language, and if a new
 10104 feature is to be added, it should be accomplished by extending the capabilities presented by the

10105 current model of the shell, rather than by inventing a new one outside the model; adding a new
 10106 dollar expansion was perceived to be the most intuitive and least destructive way to add such a
 10107 new capability.

10108 The standard requires assignment operators to be supported (as listed in the Shell and Utilities
 10109 volume of IEEE Std 1003.1-200x, Section 1.7.2, Concepts Derived from the ISO C Standard), and
 10110 since arithmetic expansions are not specified to be evaluated in a subshell environment, changes
 10111 to variables there have to be in effect after the arithmetic expansion, just as in the parameter
 10112 expansion "`{x=value}`".

10113 Note, however, that "`((x=5))`" need not be equivalent to "`(($x=5))`". If the value of
 10114 the environment variable *x* is the string "y=", the expansion of "`((x=5))`" would set *x* to 5
 10115 and output 5, but "`(($x=5))`" would output 0 if the value of the environment variable *y* is
 10116 not 5 and would output 1 if the environment variable *y* is 5. Similarly, if the value of the
 10117 environment variable is 4, the expansion of "`((x=5))`" would still set *x* to 5 and output 5,
 10118 but "`(($x=5))`" (which would be equivalent to "`((4=5))`") would yield a syntax
 10119 error.

10120 In early proposals, a form `$(expression)` was used. It was functionally equivalent to the "`(())`"
 10121 of the current text, but objections were lodged that the 1988 KornShell had already implemented
 10122 "`(())`" and there was no compelling reason to invent yet another syntax. Furthermore, the
 10123 "`$()`" syntax had a minor incompatibility involving the patterns in **case** statements.

10124 The portion of the ISO C standard arithmetic operations selected corresponds to the operations
 10125 historically supported in the KornShell. In addition to the exceptions listed in the Shell and
 10126 Utilities volume of IEEE Std 1003.1-200x, Section 2.6.4, Arithmetic Expansion, the use of the
 10127 following are explicitly outside the scope of the rules defined in the Shell and Utilities volume of
 10128 IEEE Std 1003.1-200x, Section 1.7.2.1, Arithmetic Precision and Operations:

- 10129 • The prefix operator '&' and the "[]", "->", and '.' operators.
- 10130 • Casts

10131 It was concluded that the `test` command (`(l)`) was sufficient for the majority of relational arithmetic
 10132 tests, and that tests involving complicated relational expressions within the shell are rare, yet
 10133 could still be accommodated by testing the value of "`(())`" itself. For example:

```
10134 # a complicated relational expression
10135 while [ $(( (($x + $y)/($a * $b)) < ($foo*$bar) )) -ne 0 ]
```

10136 or better yet, the rare script that has many complex relational expressions could define a
 10137 function like this:

```
10138 val() {
10139     return $((!$1))
10140 }
```

10141 and complicated tests would be less intimidating:

```
10142 while val $(( (($x + $y)/($a * $b)) < ($foo*$bar) ))
10143 do
10144     # some calculations
10145 done
```

10146 A suggestion that was not adopted was to modify `true` and `false` to take an optional argument,
 10147 and `true` would exit true only if the argument was non-zero, and `false` would exit false only if the
 10148 argument was non-zero:

```
10149 while true $(( $x > 5 && $y <= 25 ))
```

10150 There is a minor portability concern with the new syntax. The example "`$((2+2))`" could have
 10151 been intended to mean a command substitution of a utility named "2+2" in a subshell. The
 10152 standard developers considered this to be obscure and isolated to some KornShell scripts
 10153 (because "`$()`" command substitution existed previously only in the KornShell). The text on
 10154 command substitution requires that the "`$(`" and '`(`' be separate tokens if this usage is
 10155 needed.

10156 An example such as:

```
10157 echo $((echo hi);(echo there))
```

10158 should not be misinterpreted by the shell as arithmetic because attempts to balance the
 10159 parentheses pairs would indicate that they are subshells. However, as indicated by the Base
 10160 Definitions volume of IEEE Std 1003.1-200x, Section 3.112, Control Operator, a conforming
 10161 application must separate two adjacent parentheses with white space to indicate nested
 10162 subshells.

10163 The standard is intentionally silent about how a variable's numeric value in an expression is
 10164 determined from its normal "sequence of bytes" value. It could be done as a text substitution, as
 10165 a conversion like that performed by *strtol()*, or even recursive evaluation. Therefore, the only
 10166 cases for which the standard is clear are those for which both conversions produce the same
 10167 result. The cases where they give the same result are those where the sequence of bytes form a
 10168 valid integer constant. Therefore, if a variable does not contain a valid integer constant, the
 10169 behavior is unspecified.

10170 For the commands:

```
10171 x=010; echo $((x += 1))
```

10172 the output must be 9.

10173 For the commands:

```
10174 x=' 1'; echo $((x += 1))
```

10175 the results are unspecified.

10176 For the commands:

```
10177 x=1+1; echo $((x += 1))
```

10178 the results are unspecified.

10179 Although the ISO/IEC 9899:1999 standard now requires support for **long long** and allows
 10180 extended integer types with higher ranks, IEEE Std 1003.1-200x only requires arithmetic
 10181 expansions to support **signed long** integer arithmetic. Implementations are encouraged to
 10182 support signed integer values at least as large as the size of the largest file allowed on the
 10183 implementation.

10184 Implementations are also allowed to perform floating-point evaluations as long as an
 10185 application won't see different results for expressions that would not overflow **signed long**
 10186 integer expression evaluation. (This includes appropriate truncation of results to integer values.)

10187 Changes made in response to IEEE PASC Interpretation 1003.2 #208 removed the requirement
 10188 that the integer constant suffixes `l` and `L` had to be recognized. The ISO POSIX-2:1993 standard
 10189 did not require the `u`, `uL`, `uL`, `U`, `UL`, `UL`, `Lu`, `LU`, `Lu`, and `LU` suffixes since only signed integer
 10190 arithmetic was required. Since all arithmetic expressions were treated as handling **signed long**
 10191 integer types anyway, the `l` and `L` suffixes were redundant. No known scripts used them and
 10192 some historic shells did not support them. When the ISO/IEC 9899:1999 standard was used as
 10193 the basis for the description of arithmetic processing, the `ll` and `LL` suffixes and combinations
 10194 were also not required. Implementations are still free to accept any or all of these suffixes, but

10195 are not required to do so.

10196 There was also some confusion as to whether the shell was required to recognize character
 10197 constants. Syntactically, character constants were required to be recognized, but the
 10198 requirements for the handling of backslash ('\`'`) and quote ('`'`') characters (needed to specify
 10199 character constants) within an arithmetic expansion were ambiguous. Furthermore, no known
 10200 shells supported them. Changes made in response to IEEE PASC Interpretation 1003.2 #208
 10201 removed the requirement to support them (if they were indeed required before).
 10202 IEEE Std 1003.1-200x clearly does not require support for character constants.

10203 IEEE Std 1003.1-2001/Cor 2-2004, item XCU/TC2/D6/3 is applied, clarifying arithmetic
 10204 expressions.

10205 C.2.6.5 Field Splitting

10206 The operation of field splitting using *IFS*, as described in early proposals, was based on the way
 10207 the KornShell splits words, but it is incompatible with other common versions of the shell.
 10208 However, each has merit, and so a decision was made to allow both. If the *IFS* variable is unset
 10209 or is `<space><tab><newline>`, the operation is equivalent to the way the System V shell splits
 10210 words. Using characters outside the `<space><tab><newline>` set yields the KornShell behavior,
 10211 where each of the non-`<space><tab><newline>`s is significant. This behavior, which affords the
 10212 most flexibility, was taken from the way the original *awk* handled field splitting.

10213 Rule (3) can be summarized as a pseudo-ERE:

10214 $(s^*ns^* | s^+)$

10215 where *s* is an *IFS* white space character and *n* is a character in the *IFS* that is not white space.
 10216 Any string matching that ERE delimits a field, except that the *s+* form does not delimit fields at
 10217 the beginning or the end of a line. For example, if *IFS* is `<space>/<comma>/<tab>`, the string:

10218 `<space><space>red<space><space>, <space>white<space>blue`

10219 yields the three colors as the delimited fields.

10220 C.2.6.6 Pathname Expansion

10221 There is no additional rationale provided for this section.

10222 C.2.6.7 Quote Removal

10223 There is no additional rationale provided for this section.

10224 C.2.7 Redirection

10225 In the System Interfaces volume of IEEE Std 1003.1-200x, file descriptors are integers in the
 10226 range 0–(`{OPEN_MAX}`–1). The file descriptors discussed in the Shell and Utilities volume of
 10227 IEEE Std 1003.1-200x, Section 2.7, Redirection are that same set of small integers.

10228 Having multi-digit file descriptor numbers for I/O redirection can cause some obscure
 10229 compatibility problems. Specifically, scripts that depend on an example command:

10230 `echo 22>/dev/null`

10231 echoing "2" to standard error or "22" to standard output are no longer portable. However, the
 10232 file descriptor number must still be delimited from the preceding text. For example:

10233 `cat file2>foo`

10234 writes the contents of **file2**, not the contents of **file**.

10235 The "`>|`" format of output redirection was adopted from the KornShell. Along with the
 10236 *noclobber* option, *set -C*, it provides a safety feature to prevent inadvertent overwriting of

10237 existing files. (See the RATIONALE for the *pathchk* utility for why this step was taken.) The
10238 restriction on regular files is historical practice.

10239 The System V shell and the KornShell have differed historically on pathname expansion of *word*;
10240 the former never performed it, the latter only when the result was a single field (file). As a
10241 compromise, it was decided that the KornShell functionality was useful, but only as a shorthand
10242 device for interactive users. No reasonable shell script would be written with a command such
10243 as:

```
10244 cat foo > a*
```

10245 Thus, shell scripts are prohibited from doing it, while interactive users can select the shell with
10246 which they are most comfortable.

10247 The construct "2>&1" is often used to redirect standard error to the same file as standard
10248 output. Since the redirections take place beginning to end, the order of redirections is significant.
10249 For example:

```
10250 ls > foo 2>&1
```

10251 directs both standard output and standard error to file **foo**. However:

```
10252 ls 2>&1 > foo
```

10253 only directs standard output to file **foo** because standard error was duplicated as standard
10254 output before standard output was directed to file **foo**.

10255 The "<>" operator could be useful in writing an application that worked with several terminals,
10256 and occasionally wanted to start up a shell. That shell would in turn be unable to run
10257 applications that run from an ordinary controlling terminal unless it could make use of "<>"
10258 redirection. The specific example is a historical version of the pager *more*, which reads from
10259 standard error to get its commands, so standard input and standard output are both available
10260 for their usual usage. There is no way of saying the following in the shell without "<>":

```
10261 cat food | more - >/dev/tty03 2<>/dev/tty03
```

10262 Another example of "<>" is one that opens **/dev/tty** on file descriptor 3 for reading and writing:

```
10263 exec 3<> /dev/tty
```

10264 An example of creating a lock file for a critical code region:

```
10265 set -C
10266 until    2> /dev/null > lockfile
10267 do      sleep 30
10268 done
10269 set +C
10270 perform critical function
10271 rm lockfile
```

10272 Since **/dev/null** is not a regular file, no error is generated by redirecting to it in *noclobber* mode.

10273 Tilde expansion is not performed on a here-document because the data is treated as if it were
10274 enclosed in double quotes.

10275 C.2.7.1 *Redirecting Input*

10276 There is no additional rationale provided for this section.

10277 C.2.7.2 *Redirecting Output*

10278 There is no additional rationale provided for this section.

10279 C.2.7.3 *Appending Redirected Output*

10280 Note that when a file is opened (even with the `O_APPEND` flag set), the initial file offset for that
 10281 file is set to the beginning of the file. Some historic shells set the file offset to the current end-of-
 10282 file when **append** mode shell redirection was used, but this is not allowed by
 10283 IEEE Std 1003.1-200x.

10284 C.2.7.4 *Here-Document*

10285 There is no additional rationale provided for this section.

10286 C.2.7.5 *Duplicating an Input File Descriptor*

10287 There is no additional rationale provided for this section.

10288 C.2.7.6 *Duplicating an Output File Descriptor*

10289 There is no additional rationale provided for this section.

10290 C.2.7.7 *Open File Descriptors for Reading and Writing*

10291 There is no additional rationale provided for this section.

10292 **C.2.8 Exit Status and Errors**10293 C.2.8.1 *Consequences of Shell Errors*

10294 There is no additional rationale provided for this section.

10295 C.2.8.2 *Exit Status for Commands*

10296 There is a historical difference in *sh* and *ksh* non-interactive error behavior. When a command
 10297 named in a script is not found, some implementations of *sh* exit immediately, but *ksh* continues
 10298 with the next command. Thus, the Shell and Utilities volume of IEEE Std 1003.1-200x says that
 10299 the shell “may” exit in this case. This puts a small burden on the programmer, who has to test
 10300 for successful completion following a command if it is important that the next command not be
 10301 executed if the previous command was not found. If it is important for the command to have
 10302 been found, it was probably also important for it to complete successfully. The test for successful
 10303 completion would not need to change.

10304 Historically, shells have returned an exit status of $128+n$, where n represents the signal number.
 10305 Since signal numbers are not standardized, there is no portable way to determine which signal
 10306 caused the termination. Also, it is possible for a command to exit with a status in the same range
 10307 of numbers that the shell would use to report that the command was terminated by a signal.
 10308 Implementations are encouraged to choose exit values greater than 256 to indicate programs that
 10309 terminate by a signal so that the exit status cannot be confused with an exit status generated by a
 10310 normal termination.

10311 Historical shells make the distinction between “utility not found” and “utility found but cannot
 10312 execute” in their error messages. By specifying two seldomly used exit status values for these
 10313 cases, 127 and 126 respectively, this gives an application the opportunity to make use of this
 10314 distinction without having to parse an error message that would probably change from locale to
 10315 locale. The *command*, *env*, *nohup*, and *xargs* utilities in the Shell and Utilities volume of

10316 IEEE Std 1003.1-200x have also been specified to use this convention.

10317 When a command fails during word expansion or redirection, most historical implementations
 10318 exit with a status of 1. However, there was some sentiment that this value should probably be
 10319 much higher so that an application could distinguish this case from the more normal exit status
 10320 values. Thus, the language “greater than zero” was selected to allow either method to be
 10321 implemented.

10322 C.2.9 Shell Commands

10323 A description of an “empty command” was removed from an early proposal because it is only
 10324 relevant in the cases of *sh -c " "*, *system(" ")*, or an empty shell-script file (such as the
 10325 implementation of *true* on some historical systems). Since it is no longer mentioned in the Shell
 10326 and Utilities volume of IEEE Std 1003.1-200x, it falls into the silently unspecified category of
 10327 behavior where implementations can continue to operate as they have historically, but
 10328 conforming applications do not construct empty commands. (However, note that *sh* does
 10329 explicitly state an exit status for an empty string or file.) In an interactive session or a script with
 10330 other commands, extra <newline>s or semicolons, such as:

```
10331 $ false
10332 $
10333 $ echo $?
10334 1
```

10335 would not qualify as the empty command described here because they would be consumed by
 10336 other parts of the grammar.

10337 C.2.9.1 Simple Commands

10338 The enumerated list is used only when the command is actually going to be executed. For
 10339 example, in:

```
10340 true || $foo *
```

10341 no expansions are performed.

10342 The following example illustrates both how a variable assignment without a command name
 10343 affects the current execution environment, and how an assignment with a command name only
 10344 affects the execution environment of the command:

```
10345 $ x=red
10346 $ echo $x
10347 red
10348 $ export x
10349 $ sh -c 'echo $x'
10350 red
10351 $ x=blue sh -c 'echo $x'
10352 blue
10353 $ echo $x
10354 red
```

10355 This next example illustrates that redirections without a command name are still performed:

```
10356 $ ls foo
10357 ls: foo: no such file or directory
10358 $ > foo
10359 $ ls foo
10360 foo
```

10361 A command without a command name, but one that includes a command substitution, has an
 10362 exit status of the last command substitution that the shell performed. For example:

```
10363 if      x=$(command)
10364 then    ...
10365 fi
```

10366 An example of redirections without a command name being performed in a subshell shows that
 10367 the here-document does not disrupt the standard input of the **while** loop:

```
10368 IFS=:
10369 while read a b
10370 do     echo $a
10371       <<-eof
10372       Hello
10373       eof
10374 done </etc/passwd
```

10375 Following are examples of commands without command names in AND-OR lists:

```
10376 > foo || {
10377     echo "error: foo cannot be created" >&2
10378     exit 1
10379 }
```

```
10380 # set saved if /vmunix.save exists
10381 test -f /vmunix.save && saved=1
```

10382 Command substitution and redirections without command names both occur in subshells, but
 10383 they are not necessarily the same ones. For example, in:

```
10384 exec 3> file
10385 var=$(echo foo >&3) 3>&1
```

10386 it is unspecified whether **foo** is echoed to the file or to standard output.

10387 **Command Search and Execution**

10388 This description requires that the shell can execute shell scripts directly, even if the underlying
 10389 system does not support the common "#!" interpreter convention. That is, if file **foo** contains
 10390 shell commands and is executable, the following executes **foo**:

```
10391 ./foo
```

10392 The command search shown here does not match all historical implementations. A more typical
 10393 sequence has been:

- 10394 • Any built-in (special or regular)
- 10395 • Functions
- 10396 • Path search for executable files

10397 But there are problems with this sequence. Since the programmer has no idea in advance which
 10398 utilities might have been built into the shell, a function cannot be used to override portably a
 10399 utility of the same name. (For example, a function named *cd* cannot be written for many
 10400 historical systems.) Furthermore, the *PATH* variable is partially ineffective in this case, and only
 10401 a pathname with a slash can be used to ensure a specific executable file is invoked.

10402 After the *execve()* failure described, the shell normally executes the file as a shell script. Some
 10403 implementations, however, attempt to detect whether the file is actually a script and not an

10404 executable from some other architecture. The method used by the KornShell is allowed by the
10405 text that indicates non-text files may be bypassed.

10406 The sequence selected for the Shell and Utilities volume of IEEE Std 1003.1-200x acknowledges
10407 that special built-ins cannot be overridden, but gives the programmer full control over which
10408 versions of other utilities are executed. It provides a means of suppressing function lookup (via
10409 the *command* utility) for the user's own functions and ensures that any regular built-ins or
10410 functions provided by the implementation are under the control of the path search. The
10411 mechanisms for associating built-ins or functions with executable files in the path are not
10412 specified by the Shell and Utilities volume of IEEE Std 1003.1-200x, but the wording requires
10413 that if either is implemented, the application is not able to distinguish a function or built-in from
10414 an executable (other than in terms of performance, presumably). The implementation ensures
10415 that all effects specified by the Shell and Utilities volume of IEEE Std 1003.1-200x resulting from
10416 the invocation of the regular built-in or function (interaction with the environment, variables,
10417 traps, and so on) are identical to those resulting from the invocation of an executable file.

10418 IEEE Std 1003.1-2001/Cor 2-2004, item XCU/TC2/D6/4 is applied, updating the case where
10419 *execve()* fails due to an error equivalent to the [ENOEXEC] error.

10420 Examples

10421 Consider three versions of the *ls* utility:

- 10422 1. The application includes a shell function named *ls*.
- 10423 2. The user writes a utility named *ls* and puts it in **/fred/bin**.
- 10424 3. The example implementation provides *ls* as a regular shell built-in that is invoked (either
10425 by the shell or directly by *exec*) when the path search reaches the directory **/posix/bin**.

10426 If *PATH*=**/posix/bin**, various invocations yield different versions of *ls*:

10427	Invocation	Version of <i>ls</i>
10428	<i>ls</i> (from within application script)	(1) function
10429	<i>command ls</i> (from within application script)	(3) built-in
10430	<i>ls</i> (from within makefile called by application)	(3) built-in
10431	<i>system("ls")</i>	(3) built-in
10432	<i>PATH="/fred/bin:\$PATH" ls</i>	(2) user's version

10433 C.2.9.2 Pipelines

10434 Because pipeline assignment of standard input or standard output or both takes place before
10435 redirection, it can be modified by redirection. For example:

```
10436 $ command1 2>&1 | command2
```

10437 sends both the standard output and standard error of *command1* to the standard input of
10438 *command2*.

10439 The reserved word **!** allows more flexible testing using AND and OR lists.

10440 It was suggested that it would be better to return a non-zero value if any command in the
10441 pipeline terminates with non-zero status (perhaps the bitwise-inclusive OR of all return values).
10442 However, the choice of the last-specified command semantics are historical practice and would
10443 cause applications to break if changed. An example of historical behavior:

```
10444 $ sleep 5 | (exit 4)
10445 $ echo $?
10446 4
10447 $ (exit 4) | sleep 5
```

10448 \$ echo \$?
10449 0

10450 C.2.9.3 Lists

10451 The equal precedence of "&&" and "||" is historical practice. The standard developers
10452 evaluated the model used more frequently in high-level programming languages, such as C, to
10453 allow the shell logical operators to be used for complex expressions in an unambiguous way, but
10454 they could not allow historical scripts to break in the subtle way unequal precedence might
10455 cause. Some arguments were posed concerning the "{" or "(" groupings that are required
10456 historically. There are some disadvantages to these groupings:

- 10457 • The "(" can be expensive, as they spawn other processes on some implementations. This
10458 performance concern is primarily an implementation issue.
- 10459 • The "{" braces are not operators (they are reserved words) and require a trailing space
10460 after each '{', and a semicolon before each '}'. Most programmers (and certainly
10461 interactive users) have avoided braces as grouping constructs because of the problematic
10462 syntax required. Braces were not changed to operators because that would generate
10463 compatibility issues even greater than the precedence question; braces appear outside the
10464 context of a keyword in many shell scripts.

10465 IEEE PASC Interpretation 1003.2 #204 is applied, clarifying that the operators "&&" and "||"
10466 are evaluated with left associativity.

10467 Asynchronous Lists

10468 The grammar treats a construct such as:

```
10469 foo & bar & bam &
```

10470 as one "asynchronous list", but since the status of each element is tracked by the shell, the term
10471 "element of an asynchronous list" was introduced to identify just one of the **foo**, **bar**, or **bam**
10472 portions of the overall list.

10473 Unless the implementation has an internal limit, such as {CHILD_MAX}, on the retained process
10474 IDs, it would require unbounded memory for the following example:

```
10475 while true  
10476 do     foo & echo $!  
10477 done
```

10478 The treatment of the signals SIGINT and SIGQUIT with asynchronous lists is described in the
10479 Shell and Utilities volume of IEEE Std 1003.1-200x, Section 2.11, Signals and Error Handling.

10480 Since the connection of the input to the equivalent of /dev/null is considered to occur before
10481 redirections, the following script would produce no output:

```
10482 exec < /etc/passwd  
10483 cat <&0 &  
10484 wait
```

10485 **Sequential Lists**

10486 There is no additional rationale provided for this section.

10487 **AND Lists**

10488 There is no additional rationale provided for this section.

10489 **OR Lists**

10490 There is no additional rationale provided for this section.

10491 C.2.9.4 *Compound Commands*10492 **Grouping Commands**

10493 The semicolon shown in `{compound-list;}` is an example of a control operator delimiting the }
 10494 reserved word. Other delimiters are possible, as shown in the Shell and Utilities volume of
 10495 IEEE Std 1003.1-200x, Section 2.10, Shell Grammar; `<newline>` is frequently used.

10496 A proposal was made to use the **<do-done>** construct in all cases where command grouping in
 10497 the current process environment is performed, identifying it as a construct for the grouping
 10498 commands, as well as for shell functions. This was not included because the shell already has a
 10499 grouping construct for this purpose (`"{"`), and changing it would have been counter-
 10500 productive.

10501 **For Loop**

10502 The format is shown with generous usage of `<newline>`s. See the grammar in the Shell and
 10503 Utilities volume of IEEE Std 1003.1-200x, Section 2.10, Shell Grammar for a precise description of
 10504 where `<newline>`s and semicolons can be interchanged.

10505 Some historical implementations support `'{'` and `'}'` as substitutes for **do** and **done**. The
 10506 standard developers chose to omit them, even as an obsolescent feature. (Note that these
 10507 substitutes were only for the **for** command; the **while** and **until** commands could not use them
 10508 historically because they are followed by compound-lists that may contain `"{...}"` grouping
 10509 commands themselves.)

10510 The reserved word pair **do ... done** was selected rather than **do ... od** (which would have
 10511 matched the spirit of **if ... fi** and **case ... esac**) because *od* is already the name of a standard
 10512 utility.

10513 PASC Interpretation 1003.2 #169 has been applied changing the grammar.

10514 **Case Conditional Construct**

10515 An optional left parenthesis before *pattern* was added to allow numerous historical KornShell
 10516 scripts to conform. At one time, using the leading parenthesis was required if the **case** statement
 10517 was to be embedded within a `"$()"` command substitution; this is no longer the case with the
 10518 POSIX shell. Nevertheless, many historical scripts use the left parenthesis, if only because it
 10519 makes matching-parenthesis searching easier in *vi* and other editors. This is a relatively simple
 10520 implementation change that is upwards-compatible for all scripts.

10521 Consideration was given to requiring *break* inside the *compound-list* to prevent falling through to
 10522 the next pattern action list. This was rejected as being nonexistent practice. An interesting
 10523 undocumented feature of the KornShell is that using `";&"` instead of `";;"` as a terminator
 10524 causes the exact opposite behavior—the flow of control continues with the next *compound-list*.

10525 The pattern `'*'`, given as the last pattern in a **case** construct, is equivalent to the default case in
 10526 a C-language **switch** statement.

10527 The grammar shows that reserved words can be used as patterns, even if one is the first word on
 10528 a line. Obviously, the reserved word **esac** cannot be used in this manner.

10529 **If Conditional Construct**

10530 The precise format for the command syntax is described in the Shell and Utilities volume of
 10531 IEEE Std 1003.1-200x, Section 2.10, Shell Grammar.

10532 **While Loop**

10533 The precise format for the command syntax is described in the Shell and Utilities volume of
 10534 IEEE Std 1003.1-200x, Section 2.10, Shell Grammar.

10535 **Until Loop**

10536 The precise format for the command syntax is described in the Shell and Utilities volume of
 10537 IEEE Std 1003.1-200x, Section 2.10, Shell Grammar.

10538 **C.2.9.5 Function Definition Command**

10539 The description of functions in an early proposal was based on the notion that functions should
 10540 behave like miniature shell scripts; that is, except for sharing variables, most elements of an
 10541 execution environment should behave as if they were a new execution environment, and
 10542 changes to these should be local to the function. For example, traps and options should be reset
 10543 on entry to the function, and any changes to them do not affect the traps or options of the caller.
 10544 There were numerous objections to this basic idea, and the opponents asserted that functions
 10545 were intended to be a convenient mechanism for grouping common commands that were to be
 10546 executed in the current execution environment, similar to the execution of the *dot* special
 10547 built-in.

10548 It was also pointed out that the functions described in that early proposal did not provide a local
 10549 scope for everything a new shell script would, such as the current working directory, or *umask*,
 10550 but instead provided a local scope for only a few select properties. The basic argument was that
 10551 if a local scope is needed for the execution environment, the mechanism already existed: the
 10552 application can put the commands in a new shell script and call that script. All historical shells
 10553 that implemented functions, other than the KornShell, have implemented functions that operate
 10554 in the current execution environment. Because of this, traps and options have a global scope
 10555 within a shell script. Local variables within a function were considered and included in another
 10556 early proposal (controlled by the special built-in *local*), but were removed because they do not fit
 10557 the simple model developed for functions and because there was some opposition to adding yet
 10558 another new special built-in that was not part of historical practice. Implementations should
 10559 reserve the identifier *local* (as well as *typeset*, as used in the KornShell) in case this local variable
 10560 mechanism is adopted in a future version of IEEE Std 1003.1-200x.

10561 A separate issue from the execution environment of a function is the availability of that function
 10562 to child shells. A few objectors maintained that just as a variable can be shared with child shells
 10563 by exporting it, so should a function. In early proposals, the *export* command therefore had a *-f*
 10564 flag for exporting functions. Functions that were exported were to be put into the environment
 10565 as *name(=value* pairs, and upon invocation, the shell would scan the environment for these and
 10566 automatically define these functions. This facility was strongly opposed and was omitted. Some
 10567 of the arguments against exportable functions were as follows:

- 10568 • There was little historical practice. The Ninth Edition shell provided them, but there was
 10569 controversy over how well it worked.

- 10570 • There are numerous security problems associated with functions appearing in the
10571 environment of a user and overriding standard utilities or the utilities owned by the
10572 application.
- 10573 • There was controversy over requiring *make* to import functions, where it has historically
10574 used an *exec* function for many of its command line executions.
- 10575 • Functions can be big and the environment is of a limited size. (The counter-argument was
10576 that functions are no different from variables in terms of size: there can be big ones, and
10577 there can be small ones—and just as one does not export huge variables, one does not
10578 export huge functions. However, this might not apply to the average shell-function writer,
10579 who typically writes much larger functions than variables.)

10580 As far as can be determined, the functions in the Shell and Utilities volume of
10581 IEEE Std 1003.1-200x match those in System V. Earlier versions of the KornShell had two
10582 methods of defining functions:

```
10583 function fname { compound-list }
```

10584 and:

```
10585 fname() { compound-list }
```

10586 The latter used the same definition as the Shell and Utilities volume of IEEE Std 1003.1-200x, but
10587 differed in semantics, as described previously. The current edition of the KornShell aligns the
10588 latter syntax with the Shell and Utilities volume of IEEE Std 1003.1-200x and keeps the former as
10589 is.

10590 The name space for functions is limited to that of a *name* because of historical practice.
10591 Complications in defining the syntactic rules for the function definition command and in
10592 dealing with known extensions such as the "@()" usage in the KornShell prevented the name
10593 space from being widened to a *word*. Using functions to support synonyms such as the "!!"
10594 and '% ' usage in the C shell is thus disallowed to conforming applications, but acceptable as an
10595 extension. For interactive users, the aliasing facilities in the Shell and Utilities volume of
10596 IEEE Std 1003.1-200x should be adequate for this purpose. It is recognized that the name space
10597 for utilities in the file system is wider than that currently supported for functions, if the portable
10598 filename character set guidelines are ignored, but it did not seem useful to mandate extensions
10599 in systems for so little benefit to conforming applications.

10600 The "()" in the function definition command consists of two operators. Therefore, intermixing
10601 <blank>s with the *fname*, '(', and ')' is allowed, but unnecessary.

10602 An example of how a function definition can be used wherever a simple command is allowed:

```
10603 # If variable i is equal to "yes",
10604 # define function foo to be ls -l
10605 #
10606 [ "$i" = yes ] && foo() {
10607     ls -l
10608 }
```

10609 C.2.10 Shell Grammar

10610 There are several subtle aspects of this grammar where conventional usage implies rules about
10611 the grammar that in fact are not true.

10612 For *compound_list*, only the forms that end in a *separator* allow a reserved word to be recognized,
10613 so usually only a *separator* can be used where a compound list precedes a reserved word (such as
10614 **Then, Else, Do, and Rbrace**). Explicitly requiring a separator would disallow such valid (if rare)
10615 statements as:

```
10616 if (false) then (echo x) else (echo y) fi
```

10617 See the Note under special grammar rule (1).

10618 Concerning the third sentence of rule (1) (“Also, if the parser ...”):

- 10619 • This sentence applies rather narrowly: when a compound list is terminated by some clear
10620 delimiter (such as the closing **fi** of an inner **if_clause**) then it would apply; where the
10621 compound list might continue (as in after a **;**), rule (7a) (and consequently the first
10622 sentence of rule (1)) would apply. In many instances the two conditions are identical, but
10623 this part of rule (1) does not give license to treating a **WORD** as a reserved word unless it
10624 is in a place where a reserved word has to appear.
- 10625 • The statement is equivalent to requiring that when the LR(1) lookahead set contains
10626 exactly one reserved word, it must be recognized if it is present. (Here “LR(1)” refers to the
10627 theoretical concepts, not to any real parser generator.)

10628 For example, in the construct below, and when the parser is at the point marked with **^**,
10629 the only next legal token is **then** (this follows directly from the grammar rules):

```
10630 if if...fi then ... fi  
10631     ^
```

10632 At that point, the **then** must be recognized as a reserved word.

10633 (Depending on the parser generator actually used, “extra” reserved words may be in some
10634 lookahead sets. It does not really matter if they are recognized, or even if any possible
10635 reserved word is recognized in that state, because if it is recognized and is not in the
10636 (theoretical) LR(1) lookahead set, an error is ultimately detected. In the example above, if
10637 some other reserved word (for example, **while**) is also recognized, an error occurs later.

10638 This is approximately equivalent to saying that reserved words are recognized after other
10639 reserved words (because it is after a reserved word that this condition occurs), but avoids
10640 the “except for ...” list that would be required for **case**, **for**, and so on. (Reserved words
10641 are of course recognized anywhere a *simple_command* can appear, as well. Other rules take
10642 care of the special cases of non-recognition, such as rule (4) for **case** statements.)

10643 Note that the body of here-documents are handled by token recognition (see the Shell and
10644 Utilities volume of IEEE Std 1003.1-200x, Section 2.3, Token Recognition) and do not appear in
10645 the grammar directly. (However, the here-document I/O redirection operator is handled as part
10646 of the grammar.)

10647 The start symbol of the grammar (**complete_command**) represents either input from the
10648 command line or a shell script. It is repeatedly applied by the interpreter to its input and
10649 represents a single “chunk” of that input as seen by the interpreter.

10650 C.2.10.1 Shell Grammar Lexical Conventions

10651 There is no additional rationale provided for this section.

10652 C.2.10.2 Shell Grammar Rules

10653 There is no additional rationale provided for this section.

10654 C.2.11 Signals and Error Handling

10655 SD5-XCU-ERN-93 is applied, updating the first paragraph of Shell and Utilities volume of
10656 IEEE Std 1003.1-200x, Section 2.11, Signals and Error Handling.

10657 C.2.12 Shell Execution Environment

10658 Some implementations have implemented the last stage of a pipeline in the current environment
10659 so that commands such as:10660 `command | read foo`10661 set variable **foo** in the current environment. This extension is allowed, but not required;
10662 therefore, a shell programmer should consider a pipeline to be in a subshell environment, but
10663 not depend on it.10664 In early proposals, the description of execution environment failed to mention that each
10665 command in a multiple command pipeline could be in a subshell execution environment. For
10666 compatibility with some historical shells, the wording was phrased to allow an implementation
10667 to place any or all commands of a pipeline in the current environment. However, this means that
10668 a POSIX application must assume each command is in a subshell environment, but not depend
10669 on it.10670 The wording about shell scripts is meant to convey the fact that describing “trap actions” can
10671 only be understood in the context of the shell command language. Outside of this context, such
10672 as in a C-language program, signals are the operative condition, not traps.

10673 C.2.13 Pattern Matching Notation

10674 Pattern matching is a simpler concept and has a simpler syntax than REs, as the former is
10675 generally used for the manipulation of filenames, which are relatively simple collections of
10676 characters, while the latter is generally used to manipulate arbitrary text strings of potentially
10677 greater complexity. However, some of the basic concepts are the same, so this section points
10678 liberally to the detailed descriptions in the Base Definitions volume of IEEE Std 1003.1-200x,
10679 Chapter 9, Regular Expressions.

10680 C.2.13.1 Patterns Matching a Single Character

10681 Both quoting and escaping are described here because pattern matching must work in three
10682 separate circumstances:

- 10683 1. Calling directly upon the shell, such as in pathname expansion or in a
- case**
- statement. All
-
- 10684 of the following match the string or file
- abc**
- :

10685 `abc "abc" a"b" c a\bc a[b]c a["b"]c a[\b]c a["\b"]c a?c a*c`

10686 The following do not:

10687 `"a?c" a*c a[b]c`

- 10688 2. Calling a utility or function without going through a shell, as described for
- find*
- and the
-
- 10689
- fnmatch()*
- function defined in the System Interfaces volume of IEEE Std 1003.1-200x.

10690 3. Calling utilities such as *find*, *cpio*, *tar*, or *pax* through the shell command line. In this case,
10691 shell quote removal is performed before the utility sees the argument. For example, in:

```
10692 find /bin -name "e\c[\h]o" -print
```

10693 after quote removal, the backslashes are presented to *find* and it treats them as escape
10694 characters. Both precede ordinary characters, so the *c* and *h* represent themselves and *echo*
10695 would be found on many historical systems (that have it in **/bin**). To find a filename that
10696 contained shell special characters or pattern characters, both quoting and escaping are
10697 required, such as:

```
10698 pax -r ... "*a\(\?"
```

10699 to extract a filename ending with "a(?".

10700 Conforming applications are required to quote or escape the shell special characters (sometimes
10701 called metacharacters). If used without this protection, syntax errors can result or
10702 implementation extensions can be triggered. For example, the KornShell supports a series of
10703 extensions based on parentheses in patterns.

10704 The restriction on a circumflex in a bracket expression is to allow implementations that support
10705 pattern matching using the circumflex as the negation character in addition to the exclamation
10706 mark. A conforming application must use something like "[\^!]" to match either character.

10707 C.2.13.2 Patterns Matching Multiple Characters

10708 Since each asterisk matches zero or more occurrences, the patterns "a*b" and "a**b" have
10709 identical functionality.

10710 Examples

10711 a[bc] Matches the strings "ab" and "ac".

10712 a*d Matches the strings "ad", "abd", and "abcd", but not the string "abc".

10713 a*d* Matches the strings "ad", "abcd", "abcdef", "aaaad", and "adddd".

10714 *a*d Matches the strings "ad", "abcd", "efabcd", "aaaad", and "adddd".

10715 C.2.13.3 Patterns Used for Filename Expansion

10716 The caveat about a slash within a bracket expression is derived from historical practice. The
10717 pattern "a[b/c]d" does not match such pathnames as **abd** or **a/d**. On some implementations
10718 (including those conforming to the Single UNIX Specification), it matched a pathname of
10719 literally "a[b/c]d". On other systems, it produced an undefined condition (an unescaped '['
10720 used outside a bracket expression). In this version, the XSI behavior is now required.

10721 Filenames beginning with a period historically have been specially protected from view on
10722 UNIX systems. A proposal to allow an explicit period in a bracket expression to match a leading
10723 period was considered; it is allowed as an implementation extension, but a conforming
10724 application cannot make use of it. If this extension becomes popular in the future, it will be
10725 considered for a future version of the Shell and Utilities volume of IEEE Std 1003.1-200x.

10726 Historical systems have varied in their permissions requirements. To match **f*/bar** has required
10727 read permissions on the **f*** directories in the System V shell, but the Shell and Utilities volume of
10728 IEEE Std 1003.1-200x, the C shell, and KornShell require only search permissions.

10729 C.2.14 Special Built-In Utilities

10730 See the RATIONALE sections on the individual reference pages.

10731 C.3 Batch Environment Services and Utilities

10732 Scope of the Batch Environment Services and Utilities Option

10733 This section summarizes the deliberations of the IEEE P1003.15 (Batch Environment) working
10734 group in the development of the Batch Environment Services and Utilities option, which covers
10735 a set of services and utilities defining a batch processing system.

10736 This informative section contains historical information concerning the contents of the
10737 amendment and describes why features were included or discarded by the working group.

10738 History of Batch Systems

10739 The supercomputing technical committee began as a “Birds Of a Feather” (BOF) at the January
10740 1987 Usenix meeting. There was enough general interest to form a supercomputing attachment
10741 to the /usr/group working groups. Several subgroups rapidly formed. Of those subgroups, the
10742 batch group was the most ambitious. The first early meetings were spent evaluating user needs
10743 and existing batch implementations.

10744 To evaluate user needs, individuals from the supercomputing community came and presented
10745 their needs. Common requests were flexibility, interoperability, control of resources, and ease-of-
10746 use. Backward-compatibility was not an issue. The working group then evaluated some existing
10747 systems. The following different systems were evaluated:

- 10748 • PROD
- 10749 • Convex Distributed Batch
- 10750 • NQS
- 10751 • CTSS
- 10752 • MDQS from Ballistics Research Laboratory (BRL)

10753 Finally, NQS was chosen as a model because it satisfied not only the most user requirements, but
10754 because it was public domain, already implemented on a variety of hardware platforms, and
10755 network-based.

10756 Historical Implementations of Batch Systems

10757 Deferred processing of work under the control of a scheduler has been a feature of most
10758 proprietary operating systems from the earliest days of multi-user systems in order to maximize
10759 utilization of the computer.

10760 The arrival of UNIX systems proved to be a dilemma to many hardware providers and users
10761 because it did not include the sophisticated batch facilities offered by the proprietary systems.
10762 This omission was rectified in 1986 by NASA Ames Research Center who developed the
10763 Network Queuing System (NQS) as a portable UNIX application that allowed the routing and
10764 processing of batch “jobs” in a network. To encourage its usage, the product was later put into
10765 the public domain. It was promptly picked up by UNIX hardware providers, and ported and
10766 developed for their respective hardware and UNIX implementations.

10767 Many major vendors, who traditionally offer a batch-dominated environment, ported the
10768 public-domain product to their systems, customized it to support the capabilities of their
10769 systems, and added many customer-requested features.

10770 Due to the strong hardware provider and customer acceptance of NQS, it was decided to use
 10771 NQS as the basis for the POSIX Batch Environment amendment in 1987. Other batch systems
 10772 considered at the time included CTSS, MDQS (a forerunner of NQS from the Ballistics Research
 10773 Laboratory), and PROD (a Los Alamos Labs development). None were thought to have both the
 10774 functionality and acceptability of NQS.

10775 **NQS Differences from the *at* utility**

10776 The base standard *at* and *batch* utilities are not sufficient to meet the batch processing needs in a
 10777 supercomputing environment and additional functionality in the areas of resource management,
 10778 job scheduling, system management, and control of output is required.

10779 **Batch Environment Services and Utilities Option Definitions**

10780 The concept of a batch job is closely related to a session with a session leader. The main
 10781 difference is that a batch job does not have a controlling terminal. There has been much debate
 10782 over whether to use the term “request” or “job”. Job was the final choice because of the
 10783 historical use of this term in the batch environment.

10784 The current definition for job identifiers is not sufficient with the model of destinations. The
 10785 current definition is:

10786 `sequence_number.originating_host`

10787 Using the model of destination, a host may include multiple batch nodes, the location of which
 10788 is identified uniquely by a name or directory service. If the current definition is used, batch
 10789 nodes running on the same host would have to coordinate their use of sequence numbers, as
 10790 sequence numbers are assigned by the originating host. The alternative is to use the originating
 10791 batch node name instead of the originating host name.

10792 The reasons for wishing to run more than one batch system per host could be the following.

10793 A test and production batch system are maintained on a single host. This is most likely in a
 10794 development facility, but could also arise when a site is moving from one version to another. The
 10795 new batch system could be installed as a test version that is completely separate from the
 10796 production batch system, so that problems can be isolated to the test system. Requiring the batch
 10797 nodes to coordinate their use of sequence numbers creates a dependency between the two
 10798 nodes, and that defeats the purpose of running two nodes.

10799 A site has multiple departments using a single host, with different management policies. An
 10800 example of contention might be in job selection algorithms. One group might want a FIFO type
 10801 of selection, while another group wishes to use a more complex algorithm based on resource
 10802 availability. Again, requiring the batch nodes to coordinate is an unnecessary binding.

10803 The proposal eventually accepted was to replace originating host with originating batch node.
 10804 This supplies sufficient granularity to ensure unique job identifiers. If more than one batch node
 10805 is on a particular host, they each have their own unique name.

10806 The queue portion of a destination is not part of the job identifier as these are not required to be
 10807 unique between batch nodes. For instance, two batch nodes may both have queues called small,
 10808 medium, and large. It is only the batch node name that is uniquely identifiable throughout the
 10809 batch system. The queue name has no additional function in this context.

10810 Assume there are three batch nodes, each of which has its own name server. On batch node one,
 10811 there are no queues. On batch node two, there are fifty queues. On batch node three, there are
 10812 forty queues. The system administrator for batch node one does not have to configure queues,
 10813 because there are none implemented. However, if a user wishes to send a job to either batch
 10814 node two or three, the system administrator for batch node one must configure a destination

10815 that maps to the appropriate batch node and queue. If every queue is to be made accessible from
10816 batch node one, the system administrator has to configure ninety destinations.

10817 To avoid requiring this, there should be a mechanism to allow a user to separate the destination
10818 into a batch node name and a queue name. Then, an implementation that is configured to get to
10819 all the batch nodes does not need any more configuration to allow a user to get to all of the
10820 queues on all of the batch nodes. The node name is used to locate the batch node, while the
10821 queue name is sent unchanged to that batch node.

10822 The following are requirements that a destination identifier must be capable of providing:

- 10823 • The ability to direct a job to a queue in a particular batch node.
- 10824 • The ability to direct a job to a particular batch node.
- 10825 • The ability to group at a higher level than just one queue. This includes grouping similar
10826 queues across multiple batch nodes (this is a pipe queue).
- 10827 • The ability to group batch nodes. This allows a user to submit a job to a group name with
10828 no knowledge of the batch node configuration. This also provides aliasing as a special
10829 case. Aliasing is a group containing only one batch node name. The group name is the
10830 alias.

10831 In addition, the administrator has the following requirements:

- 10832 • The ability to control access to the queues.
- 10833 • The ability to control access to the batch nodes.
- 10834 • The ability to control access to groups of queues (pipe queues).
- 10835 • The ability to configure retry time intervals and durations.

10836 The requirements of the user are met by destination as explained in the following.

10837 The user has the ability to specify a queue name, which is known only to the batch node
10838 specified. There is no configuration of these queues required on the submitting node.

10839 The user has the ability to specify a batch node whose name is network-unique. The
10840 configuration required is that the batch node be defined as an application, just as other
10841 applications such as FTP are configured.

10842 Once a job reaches a queue, it can again become a user of the batch system. The batch node can
10843 choose to send the job to another batch node or queue or both. In other words, the routing is at
10844 an application level, and it is up to the batch system to choose where the job will be sent.
10845 Configuration is up to the batch node where the queue resides. This provides grouping of
10846 queues across batch nodes or within a batch node. The user submits the job to a queue, which by
10847 definition routes the job to other queues or nodes or both.

10848 A node name may be given to a naming service, which returns multiple addresses as opposed to
10849 just one. This provides grouping at a batch node level. This is a local issue, meaning that the
10850 batch node must choose only one of these addresses. The list of addresses is not sent with the
10851 job, and once the job is accepted on another node, there is no connection between the list and the
10852 job. The requirements of the administrator are met by destination as explained in the following.

10853 The control of queues is a batch system issue, and will be done using the batch administrative
10854 utilities.

10855 The control of nodes is a network issue, and will be done through whatever network facilities
10856 are available.

10857 The control of access to groups of queues (pipe queues) is covered by the control of any other

10858 queue. The fact that the job may then be sent to another destination is not relevant.

10859 The propagation of a job across more than one point-to-point connection was dropped because
10860 of its complexity and because all of the issues arising from this capability could not be resolved.
10861 It could be provided as additional functionality at some time in the future.

10862 The addition of *network* as a defined term was done to clarify the difference between a network
10863 of batch nodes as opposed to a network of hosts. A network of batch nodes is referred to as a
10864 batch system. The network refers to the actual host configuration. A single host may have
10865 multiple batch nodes.

10866 In the absence of a standard network naming convention, this option establishes its own
10867 convention for the sake of consistency and expediency. This is subject to change, should a future
10868 working group develop a standard naming convention for network pathnames.

10869 C.3.1 Batch General Concepts

10870 During the development of the Batch Environment Services and Utilities option, a number of
10871 topics were discussed at length which influenced the wording of the normative text but could
10872 not be included in the final text. The following items are some of the most significant terms and
10873 concepts of those discussed:

- 10874 • Small and Consistent Command Set

10875 Often, conventional utilities from UNIX systems have a very complicated utility syntax
10876 and usage. This can often result in confusion and errors when trying to use them. The
10877 Batch Environment Services and Utilities option utility set, on the other hand, has been
10878 paired to a small set of robust utilities with an orthogonal calling sequence.

- 10879 • Checkpoint/Restart

10880 This feature permits an already executing process to checkpoint or save its contents. Some
10881 implementations permit this at both the batch utility level (for example, checkpointing this
10882 job upon its abnormal termination) or from within the job itself via a system call. Support
10883 of checkpoint/restart is optional. A conscious, careful effort was made to make the *qsub*
10884 utility consistently refer to checkpoint/restart as optional functionality.

- 10885 • Rerunability

10886 When a user submits a job for batch processing, they can designate it “rerunnable” in that
10887 it will automatically resume execution from the start of the job if the machine on which it
10888 was executing crashes for some reason. The decision on whether the job will be rerun or
10889 not is entirely up to the submitter of the job and no decisions will be made within the batch
10890 system. A job that is rerunnable and has been submitted with the proper
10891 checkpoint/restart switch will first be checkpointed and execution begun from that point.
10892 Furthermore, use of the implementation-defined checkpoint/restart feature will not be
10893 defined in this context.

- 10894 • Error Codes

10895 All utilities exit with error status zero (0) if successful, one (1) if a user error occurred, and
10896 two (2) for an internal Batch Environment Services and Utilities option error.

- 10897 • Level of Portability

10898 Portability is specified at both the user, operator, and administrator levels. A conforming
10899 batch implementation prevents identical functionality and behavior at all these levels.
10900 Additionally, portable batch shell scripts with embedded Batch Environment Services and
10901 Utilities option utilities add an additional level of portability.

- 10902
- Resource Specification
- 10903 A small set of globally understood resources, such as memory and CPU time, is specified.
- 10904 All conforming batch implementations are able to process them in a manner consistent
- 10905 with the yet-to-be-developed resource management model. Resources not in this
- 10906 amendment set are ignored and passed along as part of the argument stream of the utility.
- 10907
- Queue Position
- 10908 Queue position is the place a job occupies in a queue. It is dependent on a variety of factors
- 10909 such as submission time and priority. Since priority may be affected by the implementation
- 10910 of fair share scheduling, the definition of queue position is implementation-defined.
- 10911
- Queue ID
- 10912 A numerical queue ID is an external requirement for purposes of accounting. The
- 10913 identification number was chosen over queue name for processing convenience.
- 10914
- Job ID
- 10915 A common notion of “jobs” is a collection of processes whose process group cannot be
- 10916 altered and is used for resource management and accounting. This concept is
- 10917 implementation-defined and, as such, has been omitted from the batch amendment.
- 10918
- Bytes *versus* Words
- 10919 Except for one case, bytes are used as the standard unit for memory size. Furthermore, the
- 10920 definition of a word varies from machine to machine. Therefore, bytes will be the default
- 10921 unit of memory size.
- 10922
- Regular Expressions
- 10923 The standard definition of regular expressions is much too broad to be used in the batch
- 10924 utility syntax. All that is needed is a simple concept of “all”; for example, delete all my jobs
- 10925 from the named queue. For this reason, regular expressions have been eliminated from the
- 10926 batch amendment.
- 10927
- Display Privacy
- 10928 How much data should be displayed locally through functions? Local policy dictates the
- 10929 amount of privacy. Library functions must be used to create and enforce local policy.
- 10930 Network and local *qstats* must reflect the policy of the server machine.
- 10931
- Remote Host Naming Convention
- 10932 It was decided that host names would be a maximum of 255 characters in length, with at
- 10933 most 15 characters being shown in displays. The 255 character limit was chosen because it
- 10934 is consistent with BSD. The 15-character limit was an arbitrary decision.
- 10935
- Network Administration
- 10936 Network administration is important, but is outside the scope of the batch amendment.
- 10937 Network administration could be done with *rsh*. However, authentication becomes two-
- 10938 sided.
- 10939
- Network Administration Philosophy
- 10940 Keep it simple. Centralized management should be possible. For example, Los Alamos
- 10941 needs a dumb set of CPUs to be managed by a central system *versus* several
- 10942 independently-managed systems as is the general case for the Batch Environment Services
- 10943 and Utilities option.

- 10944
- Operator Utility Defaults (that is, Default Host, User, Account, and so on)
- 10945 It was decided that usability would override orthogonality and syntactic consistency.
- 10946
- The Batch System Manager and Operator Distinction
- 10947 The distinction between manager and operator is that operators can only control the flow
- 10948 of jobs. A manager can alter the batch system configuration in addition to job flow. POSIX
- 10949 makes a distinction between user and system administrator but goes no further. The
- 10950 concepts of manager and operator privileges fall under local policy. The distinction
- 10951 between manager and operator is historical in batch environments, and the Batch
- 10952 Environment Services and Utilities option has continued that distinction.
- 10953
- The Batch System Administrator
- 10954 An administrator is equivalent to a batch system manager.

10955 C.3.2 Batch Services

10956 This rationale is provided as informative rather than normative text, to avoid placing

10957 requirements on implementors regarding the use of symbolic constants, but at the same time to

10958 give implementors a preferred practice for assigning values to these constants to promote

10959 interoperability.

10960 The *Checkpoint* and *Minimum_Cpu_Interval* attributes induce a variety of behavior depending

10961 upon their values. Some jobs cannot or should not be checkpointed. Other users will simply

10962 need to ensure job continuation across planned downtimes; for example, scheduled preventive

10963 maintenance. For users consuming expensive resources, or for jobs that run longer than the

10964 mean time between failures, however, periodic checkpointing may be essential. However,

10965 system administrators must be able to set minimum checkpoint intervals on a queue-by-queue

10966 basis to guard against, for example, naive users specifying interval values too small on memory-

10967 intensive jobs. Otherwise, system overhead would adversely affect performance.

10968 The use of symbolic constants, such as `NO_CHECKPOINT`, was introduced to lend a degree of

10969 formalism and portability to this option.

10970 Support for checkpointing is optional for servers. However, clients must provide for the `-c`

10971 option, since in a distributed environment the job may run on a server that does provide such

10972 support, even if the host of the client does not support the checkpoint feature.

10973 If the user does not specify the `-c` option, the default action is left unspecified by this option.

10974 Some implementations may wish to do checkpointing by default; others may wish to checkpoint

10975 only under an explicit request from the user.

10976 The *Priority* attribute has been made non-optional. All clients already had been required to

10977 support the `-p` option. The concept of prioritization is common in historical implementations.

10978 The default priority is left to the server to establish.

10979 The *Hold_Types* attribute has been modified to allow for implementation-defined hold types to

10980 be passed to a batch server.

10981 It was the intent of the IEEE P1003.15 working group to mandate the support for the

10982 *Resource_List* attribute in this option by referring to another amendment, specifically the

10983 IEEE P1003.1a draft standard. However, during the development of the IEEE P1003.1a draft

10984 standard this was excluded. As such this requirement has been removed from the normative

10985 text.

10986 The *Shell_Path* attribute has been modified to accept a list of shell paths that are associated with

10987 a host. The name of the attribute has been changed to *Shell_Path_List*.

10988 C.3.3 Common Behavior for Batch Environment Utilities

10989 This section was defined to meet the goal of a “Small and Consistent Command Set” for this
10990 option.

10991 C.4 Utilities

10992 For the utilities included in IEEE Std 1003.1-200x, see the RATIONALE sections on the
10993 individual reference pages.

10994 Exclusion of Utilities

10995 The set of utilities contained in IEEE Std 1003.1-200x is drawn from the base documents, with
10996 one addition: the *c99* utility. This section contains rationale for some of the deliberations that led
10997 to this set of utilities, and why certain utilities were excluded.

10998 Many utilities were evaluated by the standard developers; more historical utilities were
10999 excluded from the base documents than included. The following list contains many common
11000 UNIX system utilities that were not included as mandatory utilities, in the User Portability
11001 Utilities option, in the XSI option, or in one of the software development groups. It is logistically
11002 difficult for this rationale to distribute correctly the reasons for not including a utility among the
11003 various utility options. Therefore, this section covers the reasons for all utilities not included in
11004 IEEE Std 1003.1-200x.

11005 This rationale is limited to a discussion of only those utilities actively or indirectly evaluated by
11006 the standard developers of the base documents, rather than the list of all known UNIX utilities
11007 from all its variants.

11008 *adb* The intent of the various software development utilities was to assist in the
11009 installation (rather than the actual development and debugging) of applications.
11010 This utility is primarily a debugging tool. Furthermore, many useful aspects of *adb*
11011 are very hardware-specific.

11012 *as* Assemblers are hardware-specific and are included implicitly as part of the
11013 compilers in IEEE Std 1003.1-200x.

11014 *banner* The only known use of this command is as part of the *lp* printer header pages. It
11015 was decided that the format of the header is implementation-defined, so this utility
11016 is superfluous to application portability.

11017 *calendar* This reminder service program is not useful to conforming applications.

11018 *cancel* The *lp* (line printer spooling) system specified is the most basic possible and did
11019 not need this level of application control.

11020 *chroot* This is primarily of administrative use, requiring superuser privileges.

11021 *col* No utilities defined in IEEE Std 1003.1-200x produce output requiring such a filter.
11022 The *nroff* text formatter is present on many historical systems and will continue to
11023 remain as an extension; *col* is expected to be shipped by all the systems that ship
11024 *nroff*.

11025 *cpio* This has been replaced by *pax*, for reasons explained in the rationale for that utility.

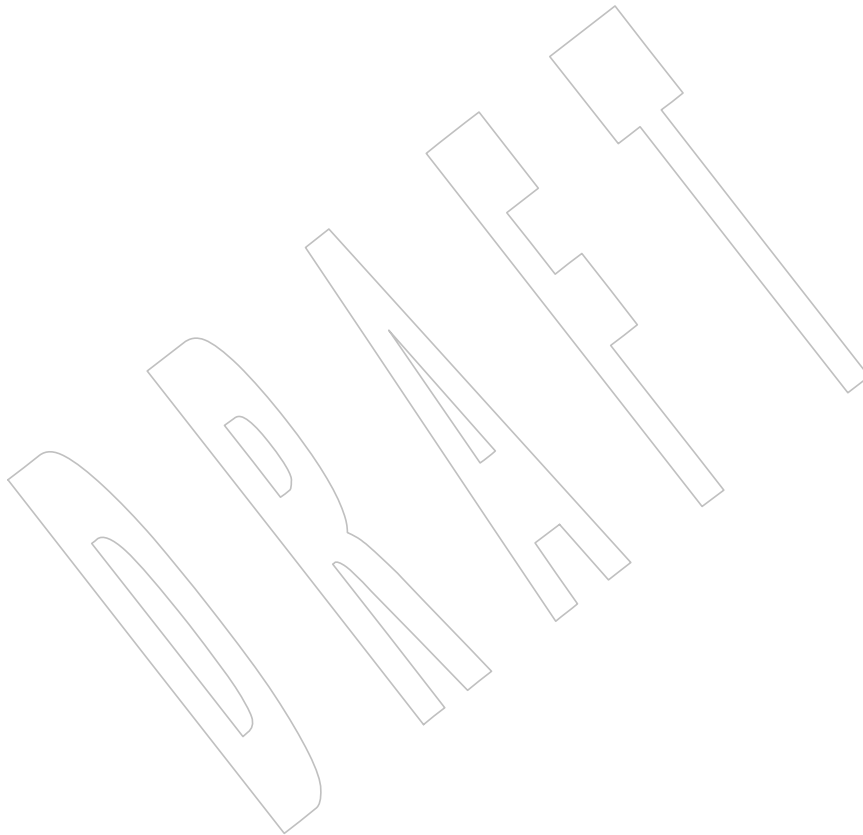
11026 *cpp* This is subsumed by *c99*.

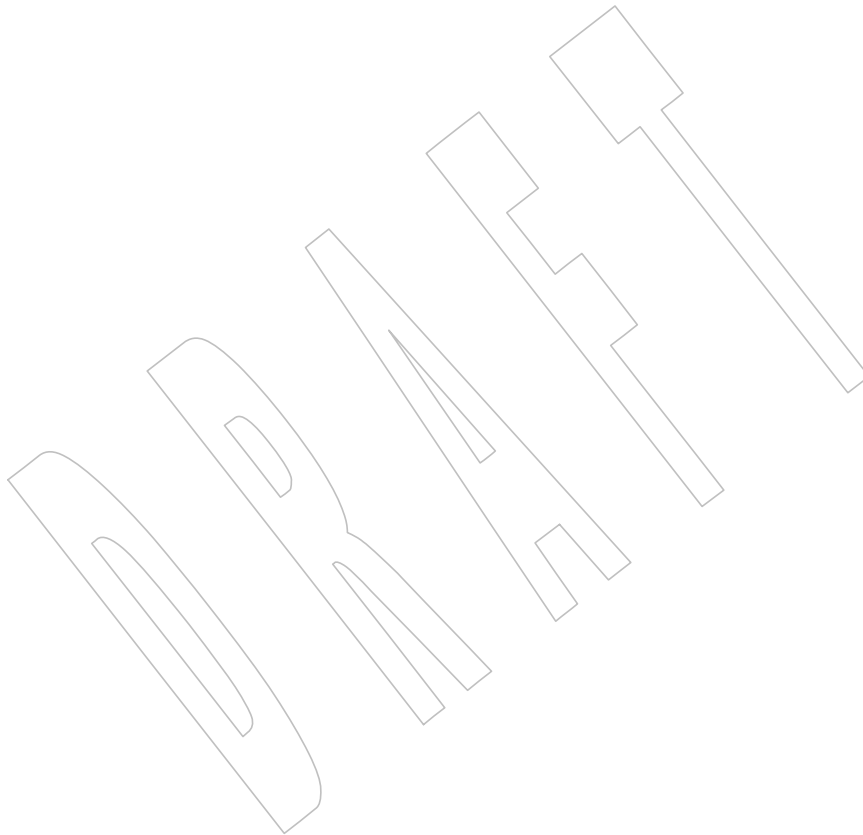
11027 *cu* This utility is terminal-oriented and is not useful from shell scripts or typical
11028 application programs.

11029	<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 1003.1-200x prescribes the interface and not the underlying mechanism used to implement it.
11030		
11031		
11032		
11033	<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.
11034		
11035		
11036	<i>dis</i>	Disassemblers are hardware-specific.
11037	<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> .
11038		
11039		
11040		
11041		
11042		
11043		
11044		
11045	<i>ld</i>	This is subsumed by <i>c99</i> .
11046	<i>line</i>	The functionality of <i>line</i> can be provided with <i>read</i> .
11047	<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.
11048		
11049		
11050		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.
11051		
11052		
11053		
11054		
11055		
11056		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).
11057		
11058		
11059		
11060		
11061		
11062		
11063	<i>login</i>	This utility is terminal-oriented and is not useful from shell scripts or typical application programs.
11064		
11065	<i>lorder</i>	This utility is an aid in creating an implementation-defined detail of object libraries that the standard developers did not feel required standardization.
11066		
11067	<i>lpstat</i>	The <i>lp</i> system specified is the most basic possible and did not need this level of application control.
11068		
11069	<i>mail</i>	This utility was omitted in favor of <i>mailx</i> because there was a considerable functionality overlap between the two.
11070		
11071	<i>mknod</i>	This was omitted in favor of <i>mkfifo</i> , as <i>mknod</i> has too many implementation-defined functions.
11072		

11073	<i>news</i>	This utility is terminal-oriented and is not useful from shell scripts or typical application programs.
11074		
11075	<i>pack</i>	This compression program was considered inferior to <i>compress</i> .
11076	<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:
11077		
11078		<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.
11079		
11080		
11081		<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 1003.1-200x for all sites using their system.
11082		
11083		
11084		
11085		
11086		
11087		<ul style="list-style-type: none"> • This is really a subject for a system administration working group or a security working group.
11088		
11089	<i>pcat</i>	This compression program was considered inferior to <i>zcat</i> .
11090	<i>pg</i>	This duplicated many of the features of the <i>more</i> pager, which was preferred by the standard developers.
11091		
11092	<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.
11093		
11094		
11095	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 option.
11096		
11097		
11098	<i>red</i>	Restricted editor. This was not considered by the standard developers because it never provided the level of security restriction required.
11099		
11100	<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.
11101		
11102		
11103	<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.
11104		
11105		
11106		
11107	<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.
11108		
11109		
11110	<i>shar</i>	Any of the numerous “shell archivers” were excluded because they did not meet the requirement of existing practice.
11111		
11112	<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.
11113		
11114		
11115	<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. This utility is primarily a debugging tool.
11116		
11117		

11118	<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.
11119		
11120		
11121		
11122	<i>su</i>	This utility is not useful from shell scripts or typical application programs. (There was also sentiment to avoid security-related utilities.)
11123		
11124	<i>sum</i>	This utility was renamed <i>cksum</i> .
11125	<i>tar</i>	This has been replaced by <i>pax</i> , for reasons explained in the rationale for that utility.
11126	<i>unpack</i>	This compression program was considered inferior to <i>uncompress</i> .
11127	<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.
11128		





11129

Rationale (Informative)

11130

Part D:

11131

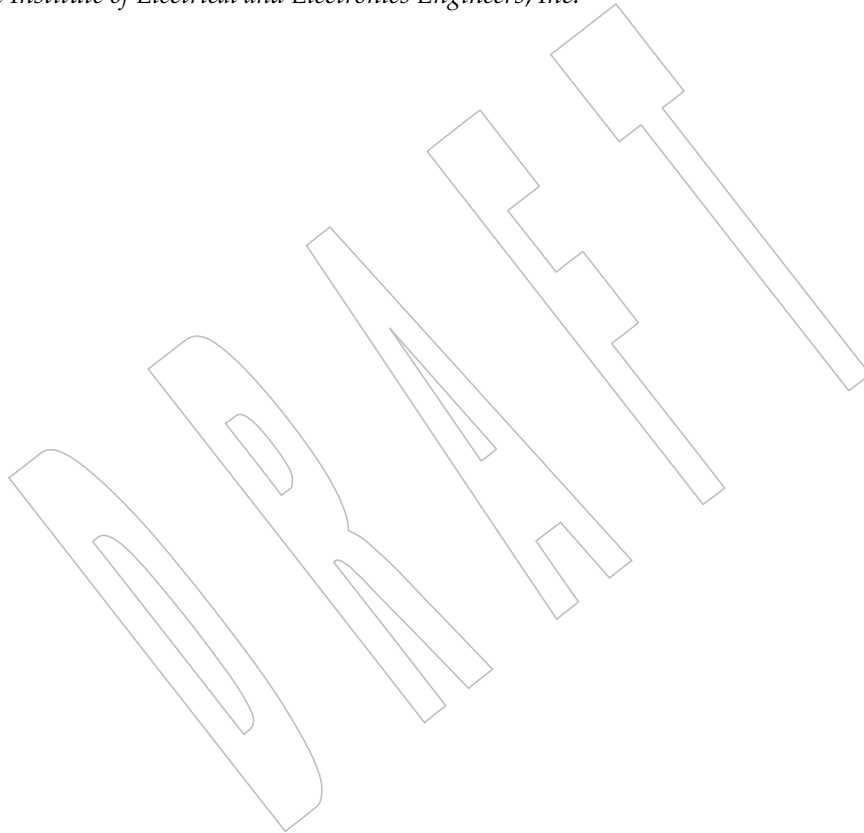
Portability Considerations

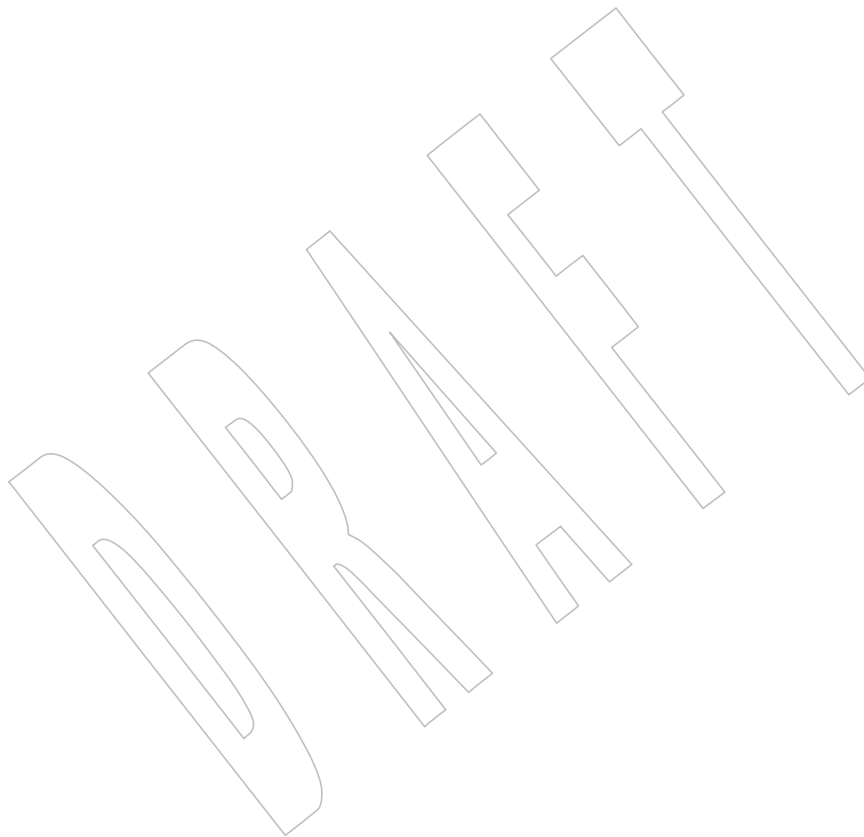
11132

The Open Group

11133

The Institute of Electrical and Electronics Engineers, Inc.





Portability Considerations (Informative)

This section contains information to satisfy various international requirements:

- [Section D.1](#) describes perceived user requirements.
- [Section D.2](#) indicates how the facilities of IEEE Std 1003.1-200x satisfy those requirements.
- [Section D.3](#) offers guidance to writers of profiles on how the configurable options, limits, and optional behavior of IEEE Std 1003.1-200x should be cited in profiles.

D.1 User Requirements

This section describes the user requirements that were perceived by the developers of IEEE Std 1003.1-200x. The primary source for these requirements was an analysis of historical practice in widespread use, as typified by the base documents listed in [Section A.1.1](#) (on page 3).

IEEE Std 1003.1-200x addresses the needs of users requiring open systems solutions for source code portability of applications. It currently addresses users requiring open systems solutions for source-code portability of applications involving multi-programming and process management (creating processes, signaling, and so on); access to files and directories in a hierarchy of file systems (opening, reading, writing, deleting files, and so on); access to asynchronous communications ports and other special devices; access to information about other users of the system; facilities supporting applications requiring bounded (realtime) response.

The following users are identified for IEEE Std 1003.1-200x:

- Those employing applications written in high-level languages, such as C, Ada, or FORTRAN.
- Users who desire conforming applications that do not necessarily require the characteristics of high-level languages (for example, the speed of execution of compiled languages or the relative security of source code intellectual property inherent in the compilation process).
- Users who desire conforming applications that can be developed quickly and can be modified readily without the use of compilers and other system components that may be unavailable on small systems or those without special application development capabilities.
- Users who interact with a system to achieve general-purpose time-sharing capabilities common to most business or government offices or academic environments: editing, filing, inter-user communications, printing, and so on.
- Users who develop applications for POSIX-conformant systems.
- Users who develop applications for UNIX systems.

An acknowledged restriction on applicable users is that they are limited to the group of individuals who are familiar with the style of interaction characteristic of historically-derived systems based on one of the UNIX operating systems (as opposed to other historical systems with different models, such as MS/DOS, Macintosh, VMS, MVS, and so on). Typical users

11173 would include program developers, engineers, or general-purpose time-sharing users.

11174 The requirements of users of IEEE Std 1003.1-200x can be summarized as a single goal:
 11175 *application source portability*. The requirements of the user are stated in terms of the requirements
 11176 of portability of applications. This in turn becomes a requirement for a standardized set of
 11177 syntax and semantics for operations commonly found on many operating systems.

11178 The following sections list the perceived requirements for application portability.

11179 **D.1.1 Configuration Interrogation**

11180 An application must be able to determine whether and how certain optional features are
 11181 provided and to identify the system upon which it is running, so that it may appropriately adapt
 11182 to its environment.

11183 Applications must have sufficient information to adapt to varying behaviors of the system.

11184 **D.1.2 Process Management**

11185 An application must be able to manage itself, either as a single process or as multiple processes.
 11186 Applications must be able to manage other processes when appropriate.

11187 Applications must be able to identify, control, create, and delete processes, and there must be
 11188 communication of information between processes and to and from the system.

11189 Applications must be able to use multiple flows of control with a process (threads) and
 11190 synchronize operations between these flows of control.

11191 **D.1.3 Access to Data**

11192 Applications must be able to operate on the data stored on the system, access it, and transmit it
 11193 to other applications. Information must have protection from unauthorized or accidental access
 11194 or modification.

11195 **D.1.4 Access to the Environment**

11196 Applications must be able to access the external environment to communicate their input and
 11197 results.

11198 **D.1.5 Access to Determinism and Performance Enhancements**

11199 Applications must have sufficient control of resource allocation to ensure the timeliness of
 11200 interactions with external objects.

11201 **D.1.6 Operating System-Dependent Profile**

11202 The capabilities of the operating system may make certain optional characteristics of the base
 11203 language in effect no longer optional, and this should be specified.

11204 **D.1.7 I/O Interaction**

11205 The interaction between the C language I/O subsystem (*stdio*) and the I/O subsystem of
 11206 IEEE Std 1003.1-200x must be specified.

- 11207 **D.1.8 Internationalization Interaction**
- 11208 The effects of the environment of IEEE Std 1003.1-200x on the internationalization facilities of the
11209 C language must be specified.
- 11210 **D.1.9 C-Language Extensions**
- 11211 Certain functions in the C language must be extended to support the additional capabilities
11212 provided by IEEE Std 1003.1-200x.
- 11213 **D.1.10 Command Language**
- 11214 Users should be able to define procedures that combine simple tools and/or applications into
11215 higher-level components that perform to the specific needs of the user. The user should be able
11216 to store, recall, use, and modify these procedures. These procedures should employ a powerful
11217 command language that is used for recurring tasks in conforming applications (scripts) in the
11218 same way that it is used interactively to accomplish one-time tasks. The language and the
11219 utilities that it uses must be consistent between systems to reduce errors and retraining.
- 11220 **D.1.11 Interactive Facilities**
- 11221 Use the system to accomplish individual tasks at an interactive terminal. The interface should be
11222 consistent, intuitive, and offer usability enhancements to increase the productivity of terminal
11223 users, reduce errors, and minimize retraining costs. Online documentation or usage assistance
11224 should be available.
- 11225 **D.1.12 Accomplish Multiple Tasks Simultaneously**
- 11226 Access applications and interactive facilities from a single terminal without requiring serial
11227 execution: switch between multiple interactive tasks; schedule one-time or periodic background
11228 work; display the status of all work in progress or scheduled; influence the priority scheduling
11229 of work, when authorized.
- 11230 **D.1.13 Complex Data Manipulation**
- 11231 Manipulate data in files in complex ways: sort, merge, compare, translate, edit, format, pattern
11232 match, select subsets (strings, columns, fields, rows, and so on). These facilities should be
11233 available to both conforming applications and interactive users.
- 11234 **D.1.14 File Hierarchy Manipulation**
- 11235 Create, delete, move/rename, copy, backup/archive, and display files and directories. These
11236 facilities should be available to both conforming applications and interactive users.
- 11237 **D.1.15 Locale Configuration**
- 11238 Customize applications and interactive sessions for the cultural and language conventions of the
11239 user. Employ a wide variety of standard character encodings. These facilities should be available
11240 to both conforming applications and interactive users.

11241 **D.1.16 Inter-User Communication**

11242 Send messages or transfer files to other users on the same system or other systems on a network.
 11243 These facilities should be available to both conforming applications and interactive users.

11244 **D.1.17 System Environment**

11245 Display information about the status of the system (activities of users and their interactive and
 11246 background work, file system utilization, system time, configuration, and presence of optional
 11247 facilities) and the environment of the user (terminal characteristics, and so on). Inform the
 11248 system operator/administrator of problems. Control access to user files and other resources.

11249 **D.1.18 Printing**

11250 Output files on a variety of output device classes, accessing devices on local or network-
 11251 connected systems. Control (or influence) the formatting, priority scheduling, and output
 11252 distribution of work. These facilities should be available to both conforming applications and
 11253 interactive users.

11254 **D.1.19 Software Development**

11255 Develop (create and manage source files, compile/interpret, debug) portable open systems
 11256 applications and package them for distribution to, and updating of, other systems.

11257 **D.2 Portability Capabilities**

11258 This section describes the significant portability capabilities of IEEE Std 1003.1-200x and
 11259 indicates how the user requirements listed in [Section D.1](#) are addressed. The capabilities are
 11260 listed in the same format as the preceding user requirements; they are summarized below:

- 11261 • Configuration Interrogation
- 11262 • Process Management
- 11263 • Access to Data
- 11264 • Access to the Environment
- 11265 • Access to Determinism and Performance Enhancements
- 11266 • Operating System-Dependent Profile
- 11267 • I/O Interaction
- 11268 • Internationalization Interaction
- 11269 • C-Language Extensions
- 11270 • Command Language
- 11271 • Interactive Facilities
- 11272 • Accomplish Multiple Tasks Simultaneously
- 11273 • Complex Data Manipulation
- 11274 • File Hierarchy Manipulation
- 11275 • Locale Configuration
- 11276 • Inter-User Communication

- 11277 • System Environment
- 11278 • Printing
- 11279 • Software Development

11280 D.2.1 Configuration Interrogation

11281 The *uname()* operation provides basic identification of the system. The *sysconf()*, *pathconf()*, and
 11282 *fpathconf()* functions and the *getconf* utility provide means to interrogate the implementation to
 11283 determine how to adapt to the environment in which it is running. These values can be either
 11284 static (indicating that all instances of the implementation have the same value) or dynamic
 11285 (indicating that different instances of the implementation have the different values, or that the
 11286 value may vary for other reasons, such as reconfiguration).

11287 Unsatisfied Requirements

11288 None directly. However, as new areas are added, there will be a need for additional capability in
 11289 this area.

11290 D.2.2 Process Management

11291 The *fork()*, *exec* family, *posix_spawn()*, and *posix_spawnp()* functions provide for the creation of
 11292 new processes or the insertion of new applications into existing processes. The *_Exit()*, *_exit()*,
 11293 *exit()*, and *abort()* functions allow for the termination of a process by itself. The *wait()*, *waitid()*,
 11294 and *waitpid()* functions allow one process to deal with the termination of another.

11295 The *times()* function allows for basic measurement of times used by a process. Various
 11296 functions, including *fstat()*, *getegid()*, *geteuid()*, *getgid()*, *getgrgid()*, *getgrnam()*, *getlogin()*,
 11297 *getpid()*, *getppid()*, *getpwnam()*, *getpwuid()*, *getuid()*, *lstat()*, and *stat()*, provide for access to the
 11298 identifiers of processes and the identifiers and names of owners of processes (and files).

11299 The various functions operating on environment variables provide for communication of
 11300 information (primarily user-configurable defaults) from a parent to child processes.

11301 The operations on the current working directory control and interrogate the directory from
 11302 which relative filename searches start. The *umask()* function controls the default protections
 11303 applied to files created by the process.

11304 The *alarm()*, *pause()*, *sleep()*, *ualarm()*, and *usleep()* operations allow the process to suspend until
 11305 a timer has expired or to be notified when a period of time has elapsed. The *time()* operation
 11306 interrogates the current time and date.

11307 The signal mechanism provides for communication of events either from other processes or
 11308 from the environment to the application, and the means for the application to control the effect
 11309 of these events. The mechanism provides for external termination of a process and for a process
 11310 to suspend until an event occurs. The mechanism also provides for a value to be associated with
 11311 an event.

11312 Job control provides a means to group processes and control them as groups, and to control their
 11313 access to the function between the user and the system (the “controlling terminal”). It also
 11314 provides the means to suspend and resume processes.

11315 The Process Scheduling option provides control of the scheduling and priority of a process.

11316 The Message Passing option provides a means for interprocess communication involving small
 11317 amounts of data.

11318 The Memory Management facilities provide control of memory resources and for the sharing of
 11319 memory. This functionality is mandatory on POSIX-conforming systems.

11320 The Threads facilities provide multiple flows of control with a process (threads),
 11321 synchronization between threads (including mutexes, barriers, and spin locks), association of
 11322 data with threads, and controlled cancellation of threads.

11323 The XSI interprocess communications functionality provide an alternate set of facilities to
 11324 manipulate semaphores, message queues, and shared memory. These are provided on XSI-
 11325 conformant systems to support conforming applications developed to run on UNIX systems.

11326 D.2.3 Access to Data

11327 The *open()*, *close()*, *fclose()*, *fopen()*, and *pipe()* functions provide for access to files and data.
 11328 Such files may be regular files, interprocess data channels (pipes), or devices. Additional types
 11329 of objects in the file system are permitted and are being contemplated for standardization.

11330 The *access()*, *chmod()*, *chown()*, *dup()*, *dup2()*, *fchmod()*, *fcntl()*, *fstat()*, *ftruncate()*, *lstat()*,
 11331 *readlink()*, *realpath()*, *stat()*, and *utime()* functions allow for control and interrogation of file and
 11332 file-related objects (including symbolic links), and their ownership, protections, and timestamps.

11333 The *fgetc()*, *fputc()*, *fread()*, *fseek()*, *fsetpos()*, *fwrite()*, *getc()*, *getchar()*, *lseek()*, *putchar()*, *putc()*,
 11334 *read()*, and *write()* functions provide for data transfer from the application to files (in all their
 11335 forms).

11336 The *closedir()*, *link()*, *mkdir()*, *opendir()*, *readdir()*, *rename()*, *rmdir()*, *rewinddir()*, and *unlink()*
 11337 functions provide for a complete set of operations on directories. Directories can arbitrarily
 11338 contain other directories, and a single file can be mentioned in more than one directory.

11339 The *faccessat()*, *openat()*, *fchmodat()*, *fchownat()*, *fstatat()*, *linkat()*, *renameat()*, *readlinkat()*,
 11340 *symlinkat()*, and *unlinkat()* functions allow for race-free and thread-safe file access. The
 11341 motivation for the introduction of these functions was as follows:

- 11342 • Interfaces taking a pathname are limited by the maximum length of a pathname
 11343 (`_SC_PATH_MAX`). The absolute path of files can far exceed this length. The alternative
 11344 solution of changing the working directory and using relative pathnames is not thread-
 11345 safe.
- 11346 • A second motivation is that files accessed outside the current working directory are subject
 11347 to attacks caused by the race condition created by changing any of the elements of the
 11348 pathnames used.
- 11349 • A third motivation is to allow application code which makes use of a virtual current
 11350 working directory for each individual thread. In the alternative model there is only one
 11351 current working directory for all threads.

11352 The file-locking mechanism provides for advisory locking (protection during transactions) of
 11353 ranges of bytes (in effect, records) in a file.

11354 The *confstr()*, *fpathconf()*, *pathconf()*, and *sysconf()* functions provide for enquiry as to the
 11355 behavior of the system where variability is permitted.

11356 The asynchronous input and output functions *aio_cancel()*, *aio_error()*, *aio_fsync()*, *aio_read()*,
 11357 *aio_return()*, *aio_suspend()*, *aio_write()*, and *lio_listio()* provide for initiation and control of
 11358 asynchronous data transfers.

11359 The Synchronized Input and Output option provides for assured commitment of data to media.

11360 D.2.4 Access to the Environment

11361 The operations and types in the Base Definitions volume of IEEE Std 1003.1-200x, Chapter 11,
 11362 General Terminal Interface are provided for access to asynchronous serial devices. The primary
 11363 intended use for these is the controlling terminal for the application (the interaction point
 11364 between the user and the system). They are general enough to be used to control any
 11365 asynchronous serial device. The functions are also general enough to be used with many other
 11366 device types as a user interface when some emulation is provided.

11367 Less detailed access is provided for other device types, but in many instances an application
 11368 need not know whether an object in the file system is a device or a regular file to operate
 11369 correctly.

11370 Unsatisfied Requirements

11371 Detailed control of common device classes, specifically magnetic tape, is not provided.

11372 D.2.5 Bounded (Realtime) Response

11373 The realtime signal functions *sigqueue()*, *sigtimedwait()*, and *sigwaitinfo()* provide queued signals
 11374 and the prioritization of the handling of signals.

11375 The SCHED_FIFO, SCHED_SPORADIC, and SCHED_RR scheduling policies provide control
 11376 over processor allocation.

11377 The semaphore functions *sem_close()*, *sem_destroy()*, *sem_getvalue()*, *sem_init()*, *sem_open()*,
 11378 *sem_post()*, *sem_timedwait()*, *sem_trywait()*, *sem_unlink()*, and *sem_wait()* provide high-
 11379 performance synchronization.

11380 The memory management functions provide memory locking for control of memory allocation,
 11381 file mapping for high performance, and shared memory for high-performance interprocess
 11382 communication. The Message Passing option provides for interprocess communication without
 11383 being dependent on shared memory.

11384 The timers functions *clock_getres()*, *clock_gettime()*, *clock_settime()*, *nanosleep()*, *timer_create()*,
 11385 *timer_delete()*, *timer_getoverrun()*, *timer_gettime()*, and *timer_settime()* provide functionality to
 11386 manipulate clocks and timers and include a high resolution function called *nanosleep()* with a
 11387 finer resolution than the *sleep()* function.

11388 The timeout functions — *pthread_mutex_timedlock()*, *pthread_rwlock_timedrdlock()*,
 11389 *pthread_rwlock_timedwrlock()*, and *sem_timedwait()* — the Typed Memory Objects option and the
 11390 Monotonic Clock option provide further facilities for applications to use to obtain predictable
 11391 bounded response.

11392 D.2.6 Operating System-Dependent Profile

11393 IEEE Std 1003.1-200x makes no distinction between text and binary files. The values of
 11394 EXIT_SUCCESS and EXIT_FAILURE are further defined.

11395 Unsatisfied Requirements

11396 None known, but the ISO C standard may contain some additional options that could be
 11397 specified.

11398 **D.2.7 I/O Interaction**

11399 IEEE Std 1003.1-200x defines how each of the ISO C standard *stdio* functions interact with the
11400 POSIX.1 operations, typically specifying the behavior in terms of POSIX.1 operations.

11401 **Unsatisfied Requirements**

11402 None.

11403 **D.2.8 Internationalization Interaction**

11404 The IEEE Std 1003.1-200x environment operations provide a means to define the environment
11405 for *setlocale()* and time functions such as *ctime()*. The *tzset()* function is provided to set time
11406 conversion information.

11407 The *nl_langinfo()* function is provided to query locale-specific cultural settings.

11408 The multiple concurrent locale functions *duplocale()*, *freelocale()*, *is*_l()*, *newlocale()*,
11409 *strcasecmp_l()*, *strcoll_l()*, *strfmon_l()*, *strncasecmp_l()*, *strxfrm_l()*, *tolower_l()*, *toupper_l()*,
11410 *towctrans_l()*, *towlower()*, *towupper()*, *uselocale()*, *wscasecmp_l()*, *wscoll_l()*, *wscncasecmp_l()*,
11411 *wcsxfrm_l()*, *wctrans_l()*, and *wctype_l()* are provide to support per-thread locale information.

11412 **Unsatisfied Requirements**

11413 None.

11414 **D.2.9 C-Language Extensions**

11415 The *setjmp()* and *longjmp()* functions are not defined to be cognizant of the signal masks defined
11416 for POSIX.1. The *sigsetjmp()* and *siglongjmp()* functions are provided to fill this gap.

11417 The *_setjmp()* and *_longjmp()* functions are provided as XSI options to support historic practice.

11418 **Unsatisfied Requirements**

11419 None.

11420 **D.2.10 Command Language**

11421 The shell command language, as described in the Shell and Utilities volume of
11422 IEEE Std 1003.1-200x, Chapter 2, Shell Command Language, is a common language useful in
11423 batch scripts, through an API to high-level languages (for the C-Language Binding option,
11424 *system()* and *popen()*) and through an interactive terminal (see the *sh* utility). The shell language
11425 has many of the characteristics of a high-level language, but it has been designed to be more
11426 suitable for user terminal entry and includes interactive debugging facilities. Through the use of
11427 pipelining, many complex commands can be constructed from combinations of data filters and
11428 other common components. Shell scripts can be created, stored, recalled, and modified by the
11429 user with simple editors.

11430 In addition to the basic shell language, the following utilities offer features that simplify and
11431 enhance programmatic access to the utilities and provide features normally found only in high-
11432 level languages: *basename*, *bc*, *command*, *dirname*, *echo*, *env*, *expr*, *false*, *printf*, *read*, *sleep*, *tee*, *test*,
11433 *time**,² *true*, *wait*, *xargs*, and all of the special built-in utilities in the Shell and Utilities volume of
11434 IEEE Std 1003.1-200x, Section 2.14, Special Built-In Utilities.

11435 2. The utilities listed with an asterisk here and later in this section are present only on systems which support the User Portability Utilities
11436 option. There may be further restrictions on the utilities offered with various configuration option combinations; see the individual utility
11437 descriptions.

11438 **Unsatisfied Requirements**

11439 None.

11440 **D.2.11 Interactive Facilities**

11441 The utilities offer a common style of command-line interface through conformance to the Utility
 11442 Syntax Guidelines (see the Base Definitions volume of IEEE Std 1003.1-200x, Section 12.2, Utility
 11443 Syntax Guidelines) and the common utility defaults (see the Shell and Utilities volume of
 11444 IEEE Std 1003.1-200x, Section 1.11, Utility Description Defaults). The *sh* utility offers an
 11445 interactive command-line history and editing facility.

11446 The following utilities can be used interactively as well as by scripts; *alias*, *fc*, *mailx*, *unalias*, and
 11447 *write*.

11448 The following utilities in the User Portability Utilities option provide for interactive use: *ex*, *more*,
 11449 and *vi*; the *man* utility offers online access to system documentation.

11450 **Unsatisfied Requirements**

11451 The command line interface to individual utilities is as intuitive and consistent as historical
 11452 practice allows. Work underway based on graphical user interfaces may be more suitable for
 11453 novice or occasional users of the system.

11454 **D.2.12 Accomplish Multiple Tasks Simultaneously**

11455 The shell command language offers background processing through the asynchronous list
 11456 command form; see the Shell and Utilities volume of IEEE Std 1003.1-200x, Section 2.9, Shell
 11457 Commands.

11458 The *nohup* utility makes background processing more robust and usable.

11459 The *kill* utility can terminate background jobs.

11460 The following utilities support periodic job scheduling, control, and display: *at*, *batch*, *crontab*,
 11461 *nice*, *ps*, and *renice*.

11462 When the User Portability Utilities option is supported, the following utilities allow
 11463 manipulation of jobs: *bg*, *fg*, and *jobs*.

11464 **Unsatisfied Requirements**

11465 Terminals with multiple windows may be more suitable for some multi-tasking interactive uses
 11466 than the job control approach in IEEE Std 1003.1-200x. See the comments on graphical user
 11467 interfaces in [Section D.2.11](#) (on page 275). The *nice* and *renice* utilities do not necessarily take
 11468 advantage of complex system scheduling algorithms that are supported by the realtime options
 11469 within IEEE Std 1003.1-200x.

11470 **D.2.13 Complex Data Manipulation**

11471 The following utilities address user requirements in this area: *asa*, *awk*, *bc*, *cmp*, *comm*, *csplit*, *cut*,
 11472 *dd*, *diff*, *ed*, *ex**, *expand*, *expr*, *find*, *fold*, *grep*, *head*, *join*, *od*, *paste*, *pr*, *printf*, *sed*, *sort*, *split*, *tabs*, *tail*, *tr*,
 11473 *unexpand*, *uniq*, *uudecode*, *uuencode*, and *wc*.

11474 **Unsatisfied Requirements**

11475 Sophisticated text formatting utilities, such as *troff* or *TeX*, are not included. Standards work in
 11476 the area of SGML may satisfy this.

11477 **D.2.14 File Hierarchy Manipulation**

11478 The following utilities address user requirements in this area: *basename*, *cd*, *chgrp*, *chmod*, *chown*,
 11479 *cksum*, *cp*, *dd*, *df*, *diff*, *dirname*, *du*, *find*, *ls*, *ln*, *mkdir*, *mkfifo*, *mv*, *patch*, *pathchk*, *pax*, *pwd*, *rm*, *rmdir*,
 11480 *test*, and *touch*.

11481 **Unsatisfied Requirements**

11482 Some graphical user interfaces offer more intuitive file manager components that allow file
 11483 manipulation through the use of icons for novice users.

11484 **D.2.15 Locale Configuration**

11485 The standard utilities are affected by the various *LC_* variables to achieve locale-dependent
 11486 operation: character classification, collation sequences, regular expressions and shell pattern
 11487 matching, date and time formats, numeric formatting, and monetary formatting. When the
 11488 *POSIX2_LOCALEDEF* option is supported, applications can provide their own locale definition
 11489 files.

11490 The following utilities address user requirements in this area: *date*, *ed*, *ex**, *find*, *grep*, *locale*,
 11491 *localedef*, *more**, *sed*, *sh*, *sort*, *tr*, *uniq*, and *vi**.

11492 The *iconv()*, *iconv_close()*, and *iconv_open()* functions are available to allow an application to
 11493 convert character data between supported character sets.

11494 The *genccat* utility and the *catopen()*, *catclose()*, and *catgets()* functions provide for message
 11495 catalog manipulation.

11496 **Unsatisfied Requirements**

11497 Some aspects of multi-byte character and state-encoded character encodings have not yet been
 11498 addressed. The C-language functions, such as *getopt()*, are generally limited to single-byte
 11499 characters. The effect of the *LC_MESSAGES* variable on message formats is only suggested at
 11500 this time.

11501 **D.2.16 Inter-User Communication**

11502 The following utilities address user requirements in this area: *cksum*, *mailx*, *mesg*, *patch*, *pax*, *talk*,
 11503 *uudecode*, *uuencode*, *who*, and *write*.

11504 The historical UUCP utilities are included as a separate UUCP Utilities option.

11505 **Unsatisfied Requirements**

11506 None.

11507 D.2.17 System Environment

11508 The following utilities address user requirements in this area: *chgrp*, *chmod*, *chown*, *df*, *du*, *env*,
11509 *getconf*, *id*, *logger*, *logname*, *mesg*, *newgrp*, *ps*, *stty*, *tput*, *tty*, *umask*, *uname*, and *who*.

11510 The *closelog()*, *openlog()*, *setlogmask()*, and *syslog()* functions provide system logging facilities on
11511 XSI-conformant systems; these are analogous to the *logger* utility.

11512 Unsatisfied Requirements

11513 None.

11514 D.2.18 Printing

11515 The following utilities address user requirements in this area: *pr* and *lp*.

11516 Unsatisfied Requirements

11517 There are no features to control the formatting or scheduling of the print jobs.

11518 D.2.19 Software Development

11519 The following utilities address user requirements in this area: *ar*, *asa*, *awk*, *c99*, *ctags*, *fort77*,
11520 *getconf*, *getopts*, *lex*, *localedef*, *make*, *nm*, *od*, *patch*, *pax*, *strings*, *strip*, *time*, and *yacc*.

11521 The *system()*, *popen()*, *pclose()*, *regcomp()*, *regexec()*, *regerror()*, *regfree()*, *fnmatch()*, *getopt()*,
11522 *glob()*, *globfree()*, *wordexp()*, and *wordfree()* functions allow C-language programmers to access
11523 some of the interfaces used by the utilities, such as argument processing, regular expressions,
11524 and pattern matching.

11525 The SCCS source-code control system utilities are available on systems supporting the XSI
11526 Development option.

11527 Unsatisfied Requirements

11528 There are no language-specific development tools related to languages other than C and
11529 FORTRAN. The C tools are more complete and varied than the FORTRAN tools. There is no
11530 data dictionary or other CASE-like development tools.

11531 D.2.20 Future Growth

11532 It is arguable whether or not all functionality to support applications is potentially within the
11533 scope of IEEE Std 1003.1-200x. As a simple matter of practicality, it cannot be. Areas such as
11534 graphics, application domain-specific functionality, windowing, and so on, should be in unique
11535 standards. As such, they are properly “Unsatisfied Requirements” in terms of providing fully
11536 conforming applications, but ones which are outside the scope of IEEE Std 1003.1-200x.

11537 However, as the standards evolve, certain functionality once considered “exotic” enough to be
11538 part of a separate standard become common enough to be included in a core standard such as
11539 this. Realtime and networking, for example, have both moved from separate standards (with
11540 much difficult cross-referencing) into IEEE Std 1003.1 over time, and although no specific areas
11541 have been identified for inclusion in future revisions, such inclusions seem likely.

11542 D.3 Profiling Considerations

11543 This section offers guidance to writers of profiles on how the configurable options, limits, and
 11544 optional behavior of IEEE Std 1003.1-200x should be cited in profiles. Profile writers should
 11545 consult the general guidance in POSIX.0 when writing POSIX Standardized Profiles.

11546 The information in this section is an inclusive list of features that should be considered by profile
 11547 writers. Subsetting of IEEE Std 1003.1-200x should follow the Base Definitions volume of
 11548 IEEE Std 1003.1-200x, Section 2.1.5.1, Subprofiling Considerations. A set of profiling options is
 11549 described in [Appendix E](#) (on page 291).

11550 D.3.1 Configuration Options

11551 There are two set of options suggested by IEEE Std 1003.1-200x: those for POSIX-conforming
 11552 systems and those for X/Open System Interface (XSI) conformance. The requirements for XSI
 11553 conformance are documented in the Base Definitions volume of IEEE Std 1003.1-200x and not
 11554 discussed further here, as they superset the POSIX conformance requirements.

11555 D.3.2 Configuration Options (Shell and Utilities)

11556 There are three broad optional configurations for the Shell and Utilities volume of
 11557 IEEE Std 1003.1-200x: basic execution system, development system, and user portability
 11558 interactive system. The options to support these, and other minor configuration options, are
 11559 listed in the Base Definitions volume of IEEE Std 1003.1-200x, Chapter 2, Conformance. Profile
 11560 writers should consult the following list and the comments concerning user requirements
 11561 addressed by various components in [Section D.2](#) (on page 270).

11562 POSIX2_UPE

11563 The system supports the User Portability Utilities option.

11564 This option is a requirement for a user portability interactive system. It is required
 11565 frequently except for those systems, such as embedded realtime or dedicated application
 11566 systems, that support little or no interactive time-sharing work by users or operators. XSI-
 11567 conformant systems support this option.

11568 POSIX2_SW_DEV

11569 The system supports the Software Development Utilities option.

11570 This option is required by many systems, even those in which actual software development
 11571 does not occur. The *make* utility, in particular, is required by many application software
 11572 packages as they are installed onto the system. If POSIX2_C_DEV is supported,
 11573 POSIX2_SW_DEV is almost a mandatory requirement because of *ar* and *make*.

11574 POSIX2_C_BIND

11575 The system supports the C-Language Bindings option.

11576 This option is required on some implementations developing complex C applications or on
 11577 any system installing C applications in source form that require the functions in this option.
 11578 The *system()* and *popen()* functions, in particular, are widely used by applications; the
 11579 others are rather more specialized.

11580 POSIX2_C_DEV

11581 The system supports the C-Language Development Utilities option.

11582 This option is required by many systems, even those in which actual C-language software
 11583 development does not occur. The *c99* utility, in particular, is required by many application
 11584 software packages as they are installed onto the system. The *lex* and *yacc* utilities are used
 11585 less frequently.

- 11586 POSIX2_FORT_DEV
 11587 The system supports the FORTRAN Development Utilities option
 11588 As with C, this option is needed on any system developing or installing FORTRAN
 11589 applications in source form.
- 11590 POSIX2_FORT_RUN
 11591 The system supports the FORTRAN Runtime Utilities option.
 11592 This option is required for some FORTRAN applications that need the *asa* utility to convert
 11593 Hollerith printing statement output. It is unknown how frequently this occurs.
- 11594 POSIX2_LOCALEDEF
 11595 The system supports the creation of locales.
 11596 This option is needed if applications require their own customized locale definitions to
 11597 operate. It is presently unknown whether many applications are dependent on this.
 11598 However, the option is virtually mandatory for systems in which internationalized
 11599 applications are developed.
 11600 XSI-conformant systems support this option.
- 11601 POSIX2_PBS
 11602 The system supports the Batch Environment Services and Utilities option.
- 11603 POSIX2_PBS_ACCOUNTING
 11604 The system supports the optional feature of accounting within the Batch Environment
 11605 Services and Utilities option. It will be required in servers that implement the optional
 11606 feature of accounting.
- 11607 POSIX2_PBS_CHECKPOINT
 11608 The system supports the optional feature of checkpoint/restart within the Batch
 11609 Environment Services and Utilities option.
- 11610 POSIX2_PBS_LOCATE
 11611 The system supports the optional feature of locating batch jobs within the Batch
 11612 Environment Services and Utilities option.
- 11613 POSIX2_PBS_MESSAGE
 11614 The system supports the optional feature of sending messages to batch jobs within the Batch
 11615 Environment Services and Utilities option.
- 11616 POSIX2_PBS_TRACK
 11617 The system supports the optional feature of tracking batch jobs within the Batch
 11618 Environment Services and Utilities option.
- 11619 POSIX2_CHAR_TERM
 11620 The system supports at least one terminal type capable of all operations described in
 11621 IEEE Std 1003.1-200x.
 11622 On systems with POSIX2_UPE, this option is almost always required. It was developed
 11623 solely to allow certain specialized vendors and user applications to bypass the requirement
 11624 for general-purpose asynchronous terminal support. For example, an application and
 11625 system that was suitable for block-mode terminals, such as IBM 3270s, would not need this
 11626 option.
 11627 XSI-conformant systems support this option.

11628 D.3.3 Configurable Limits

11629 Very few of the limits need to be increased for profiles. No profile can cite lower values.

11630 {POSIX2_BC_BASE_MAX}

11631 {POSIX2_BC_DIM_MAX}

11632 {POSIX2_BC_SCALE_MAX}

11633 {POSIX2_BC_STRING_MAX}

11634 No increase is anticipated for any of these *bc* values, except for very specialized applications
11635 involving huge numbers.

11636 {POSIX2_COLL_WEIGHTS_MAX}

11637 Some natural languages with complex collation requirements require an increase from the
11638 default 2 to 4; no higher numbers are anticipated.

11639 {POSIX2_EXPR_NEST_MAX}

11640 No increase is anticipated.

11641 {POSIX2_LINE_MAX}

11642 This number is much larger than most historical applications have been able to use. At some
11643 future time, applications may be rewritten to take advantage of even larger values.

11644 {POSIX2_RE_DUP_MAX}

11645 No increase is anticipated.

11646 {POSIX2_VERSION}

11647 This is actually not a limit, but a standard version stamp. Generally, a profile should specify
11648 the Shell and Utilities volume of IEEE Std 1003.1-200x, Chapter 2, Shell Command
11649 Language by name in the normative references section, not this value.

11650 D.3.4 Configuration Options (System Interfaces)

11651 {NGROUPS_MAX}

11652 A non-zero value indicates that the implementation supports supplementary groups.

11653 This option is needed where there is a large amount of shared use of files, but where a
11654 certain amount of protection is needed. Many profiles³ are known to require this option; it
11655 should only be required if needed, but it should never be prohibited.

11656 _POSIX_ADVISORY_INFO

11657 The system provides advisory information for file management.

11658 This option allows the application to specify advisory information that can be used to
11659 achieve better or even deterministic response time in file manager or input and output
11660 operations.

11661 _POSIX_ASYNCHRONOUS_IO

11662 Support for asynchronous input and output is mandatory in IEEE Std 1003.1-200x.

11663 _POSIX_BARRIERS

11664 Support for barrier synchronization is mandatory in IEEE Std 1003.1-200x.

11665 This facility allows efficient synchronization of multiple parallel threads in multi-processor
11666 systems in which the operation is supported in part by the hardware architecture.

11667 3. There are no formally approved profiles of IEEE Std 1003.1-200x at the time of publication; the reference here is to various profiles
11668 generated by private bodies or governments.

- 11669 `_POSIX_CHOWN_RESTRICTED`
 11670 The system restricts the right to “give away” files to other users. It is mandatory that an
 11671 implementation be able to support this facility in IEEE Std 1003.1-200x; however, it is
 11672 recognized that implementations need not enable the functionality by default.
- 11673 Some applications expect that they can change the ownership of files in this way. It is
 11674 provided where either security or system account requirements cause this ability to be a
 11675 problem. It is also known to be specified in many profiles.
- 11676 `_POSIX_CLOCK_SELECTION`
 11677 Support for clock selection is mandatory in IEEE Std 1003.1-200x.
- 11678 This facility allows applications to request a high resolution sleep in order to suspend a
 11679 thread during a relative time interval, or until an absolute time value, using the desired
 11680 clock. It also allows the application to select the clock used in a `pthread_cond_timedwait()`
 11681 function call.
- 11682 `_POSIX_CPUTIME`
 11683 The system supports the Process CPU-Time Clocks option.
- 11684 This option allows applications to use a new clock that measures the execution times of
 11685 processes or threads, and the possibility to create timers based upon these clocks, for
 11686 runtime detection (and treatment) of execution time overruns.
- 11687 `_POSIX_FSYNC`
 11688 The system supports file synchronization requests.
- 11689 This option was created to support historical systems that did not provide the feature.
 11690 Applications that are expecting guaranteed completion of their input and output operations
 11691 should require the `_POSIX_SYNC_IO` option. This option should never be prohibited.
- 11692 XSI-conformant systems support this option.
- 11693 `_POSIX_IPV6`
 11694 The system supports facilities related to Internet Protocol Version 6 (IPv6).
- 11695 This option was created to allow systems to transition to IPv6.
- 11696 `_POSIX_JOB_CONTROL`
 11697 Support for job control is mandatory in IEEE Std 1003.1-200x.
- 11698 Most applications that use it can run when it is not present, although with a degraded level
 11699 of user convenience.
- 11700 `_POSIX_MAPPED_FILES`
 11701 Support for memory mapped files is mandatory in IEEE Std 1003.1-200x.
- 11702 This facility provides for the mapping of regular files into the process address space.
- 11703 Both this facility and the Shared Memory Objects option provide shared access to memory
 11704 objects in the process address space. The `mmap()` and `munmap()` functions provide the
 11705 functionality of existing practice for mapping regular files. This functionality was deemed
 11706 unnecessary, if not inappropriate, for embedded systems applications and is expected to be
 11707 optional in subprofiles.
- 11708 `_POSIX_MEMLOCK`
 11709 The system supports the locking of the address space.
- 11710 This option was created to support historical systems that did not provide the feature. It
 11711 should only be required if needed, but it should never be prohibited.

- 11712 `_POSIX_MEMLOCK_RANGE`
 11713 The system supports the locking of specific ranges of the address space.
- 11714 For applications that have well-defined sections that need to be locked and others that do
 11715 not, IEEE Std 1003.1-200x supports an optional set of functions to lock or unlock a range of
 11716 process addresses. The following are two reasons for having a means to lock down a
 11717 specific range:
- 11718 1. An asynchronous event handler function that must respond to external events in a
 11719 deterministic manner such that page faults cannot be tolerated
 - 11720 2. An input/output “buffer” area that is the target for direct-to-process I/O, and the
 11721 overhead of implicit locking and unlocking for each I/O call cannot be tolerated
- 11722 It should only be required if needed, but it should never be prohibited.
- 11723 `_POSIX_MEMORY_PROTECTION`
 11724 Support for memory protection is mandatory in IEEE Std 1003.1-200x.
- 11725 The provision of this facility typically imposes additional hardware requirements.
- 11726 `_POSIX_PRIORITIZED_IO`
 11727 The system provides prioritization for input and output operations.
- 11728 The use of this option may interfere with the ability of the system to optimize input and
 11729 output throughput. It should only be required if needed, but it should never be prohibited.
- 11730 `_POSIX_MESSAGE_PASSING`
 11731 The system supports the passing of messages between processes.
- 11732 This option was created to support historical systems that did not provide the feature. The
 11733 functionality adds a high-performance XSI interprocess communication facility for local
 11734 communication. It should only be required if needed, but it should never be prohibited.
- 11735 `_POSIX_MONOTONIC_CLOCK`
 11736 The system supports the Monotonic Clock option.
- 11737 This option allows realtime applications to rely on a monotonically increasing clock that
 11738 does not jump backwards, and whose value does not change except for the regular ticking
 11739 of the clock.
- 11740 `_POSIX_PRIORITY_SCHEDULING`
 11741 The system provides priority-based process scheduling.
- 11742 Support of this option provides predictable scheduling behavior, allowing applications to
 11743 determine the order in which processes that are ready to run are granted access to a
 11744 processor. It should only be required if needed, but it should never be prohibited.
- 11745 `_POSIX_REALTIME_SIGNALS`
 11746 Support for realtime signals is mandatory in IEEE Std 1003.1-200x.
- 11747 This facility provides prioritized, queued signals with associated data values.
- 11748 `_POSIX_REGEX`
 11749 Support for regular expression facilities is mandatory in IEEE Std 1003.1-200x.
- 11750 `_POSIX_SAVED_IDS`
 11751 Support for this feature is mandatory in IEEE Std 1003.1-200x.
- 11752 Certain classes of applications rely on it for proper operation, and there is no alternative
 11753 short of giving the application root privileges on most implementations that did not provide
 11754 `_POSIX_SAVED_IDS`.

- 11755 _posix_semaphores
11756 Support for counting semaphores is mandatory in IEEE Std 1003.1-200x.
- 11757 _posix_shared_memory_objects
11758 The system supports the mapping of shared memory objects into the process address space.

11759 Both this option and the Memory Mapped Files option provide shared access to memory
11760 objects in the process address space. The functions defined under this option provide the
11761 functionality of existing practice for shared memory objects. This functionality was deemed
11762 appropriate for embedded systems applications and, hence, is provided under this option.
11763 It should only be required if needed, but it should never be prohibited.
- 11764 _posix_shell
11765 Support for the *sh* utility command line interpreter is mandatory in IEEE Std 1003.1-200x.
- 11766 _posix_spawn
11767 The system supports the spawn option.

11768 This option provides applications with an efficient mechanism to spawn execution of a new
11769 process.
- 11770 _posix_spinlocks
11771 Support for spin locks is mandatory in IEEE Std 1003.1-200x.

11772 This facility provides a simple and efficient synchronization mechanism for threads
11773 executing in multi-processor systems.
- 11774 _posix_sporadic_server
11775 The system supports the sporadic server scheduling policy.

11776 This option provides applications with a new scheduling policy for scheduling aperiodic
11777 processes or threads in hard realtime applications.
- 11778 _posix_synchronized_io
11779 The system supports guaranteed file synchronization.

11780 This option was created to support historical systems that did not provide the feature.
11781 Applications that are expecting guaranteed completion of their input and output operations
11782 should require this option, rather than the File Synchronization option. It should only be
11783 required if needed, but it should never be prohibited.
- 11784 _posix_threads
11785 Support for multiple threads of control within a single process is mandatory in
11786 IEEE Std 1003.1-200x.
- 11787 _posix_thread_attr_stackaddr
11788 The system supports specification of the stack address for a created thread.

11789 Applications may take advantage of support of this option for performance benefits, but
11790 dependence on this feature should be minimized. This option should never be prohibited.
- 11791 XSI-conformant systems support this option.
- 11792 _posix_thread_attr_stacksize
11793 The system supports specification of the stack size for a created thread.

11794 Applications may require this option in order to ensure proper execution, but such usage
11795 limits portability and dependence on this feature should be minimized. It should only be
11796 required if needed, but it should never be prohibited.
- 11797 XSI-conformant systems support this option.

- 11798 _POSIX_THREAD_PRIORITY_SCHEDULING
11799 The system provides priority-based thread scheduling.
- 11800 Support of this option provides predictable scheduling behavior, allowing applications to
11801 determine the order in which threads that are ready to run are granted access to a processor.
11802 It should only be required if needed, but it should never be prohibited.
- 11803 _POSIX_THREAD_PRIO_INHERIT
11804 The system provides mutual-exclusion operations with priority inheritance.
- 11805 Support of this option provides predictable scheduling behavior, allowing applications to
11806 determine the order in which threads that are ready to run are granted access to a processor.
11807 It should only be required if needed, but it should never be prohibited.
- 11808 _POSIX_THREAD_PRIO_PROTECT
11809 The system supports a priority ceiling emulation protocol for mutual-exclusion operations.
- 11810 Support of this option provides predictable scheduling behavior, allowing applications to
11811 determine the order in which threads that are ready to run are granted access to a processor.
11812 It should only be required if needed, but it should never be prohibited.
- 11813 _POSIX_THREAD_PROCESS_SHARED
11814 The system provides shared access among multiple processes to synchronization objects.
- 11815 This option was created to support historical systems that did not provide the feature. It
11816 should only be required if needed, but it should never be prohibited.
- 11817 XSI-conformant systems support this option.
- 11818 _POSIX_THREAD_SAFE_FUNCTIONS
11819 Support for thread-safe functions is mandatory in IEEE Std 1003.1-200x.
- 11820 _POSIX_THREAD_SPORADIC_SERVER
11821 The system supports the thread sporadic server scheduling policy.
- 11822 Support for this option provides applications with a new scheduling policy for scheduling
11823 aperiodic threads in hard realtime applications.
- 11824 _POSIX_TIMEOUTS
11825 Support for timeouts for some blocking services is mandatory in IEEE Std 1003.1-200x.
- 11826 _POSIX_TIMERS
11827 Support for higher resolution clocks with multiple timers per process is mandatory in
11828 IEEE Std 1003.1-200x.
- 11829 This facility is appropriate for applications requiring higher resolution timestamps or
11830 needing to control the timing of multiple activities.
- 11831 _POSIX_TRACE
11832 The system supports the Trace option.
- 11833 This option was created to allow applications to perform tracing.
- 11834 _POSIX_TRACE_EVENT_FILTER
11835 The system supports the Trace Event Filter option.
- 11836 This option is dependent on support of the Trace option.
- 11837 _POSIX_TRACE_INHERIT
11838 The system supports the Trace Inherit option.
- 11839 This option is dependent on support of the Trace option.

- 11840 `_POSIX_TRACE_LOG`
 11841 The system supports the Trace Log option.
 11842 This option is dependent on support of the Trace option.
- 11843 `_POSIX_TYPED_MEMORY_OBJECTS`
 11844 The system supports the Typed Memory Objects option.
 11845 This option was created to allow realtime applications to access different kinds of physical
 11846 memory, and allow processes in these applications to share portions of this memory.

11847 D.3.5 Configurable Limits

11848 In general, the configurable limits in the `<limits.h>` header defined in the Base Definitions
 11849 volume of IEEE Std 1003.1-200x have been set to minimal values; many applications or
 11850 implementations may require larger values. No profile can cite lower values.

11851 `{AIO_LISTIO_MAX}`
 11852 The current minimum is likely to be inadequate for most applications. It is expected that
 11853 this value will be increased by profiles requiring support for list input and output
 11854 operations.

11855 `{AIO_MAX}`
 11856 The current minimum is likely to be inadequate for most applications. It is expected that
 11857 this value will be increased by profiles requiring support for asynchronous input and
 11858 output operations.

11859 `{AIO_PRIO_DELTA_MAX}`
 11860 The functionality associated with this limit is needed only by sophisticated applications. It
 11861 is not expected that this limit would need to be increased under a general-purpose profile.

11862 `{ARG_MAX}`
 11863 The current minimum is likely to need to be increased for profiles, particularly as larger
 11864 amounts of information are passed through the environment. Many implementations are
 11865 believed to support larger values.

11866 `{CHILD_MAX}`
 11867 The current minimum is suitable only for systems where a single user is not running
 11868 applications in parallel. It is significantly too low for any system also requiring windows,
 11869 and if `_POSIX_JOB_CONTROL` is specified, it should be raised.

11870 `{CLOCKRES_MIN}`
 11871 It is expected that profiles will require a finer granularity clock, perhaps as fine as 1 μ s,
 11872 represented by a value of 1 000 for this limit.

11873 `{DELAYTIMER_MAX}`
 11874 It is believed that most implementations will provide larger values.

11875 `{LINK_MAX}`
 11876 For most applications and usage, the current minimum is adequate. Many implementations
 11877 have a much larger value, but this should not be used as a basis for raising the value unless
 11878 the applications to be used require it.

11879 `{LOGIN_NAME_MAX}`
 11880 This is not actually a limit, but an implementation parameter. No profile should impose a
 11881 requirement on this value.

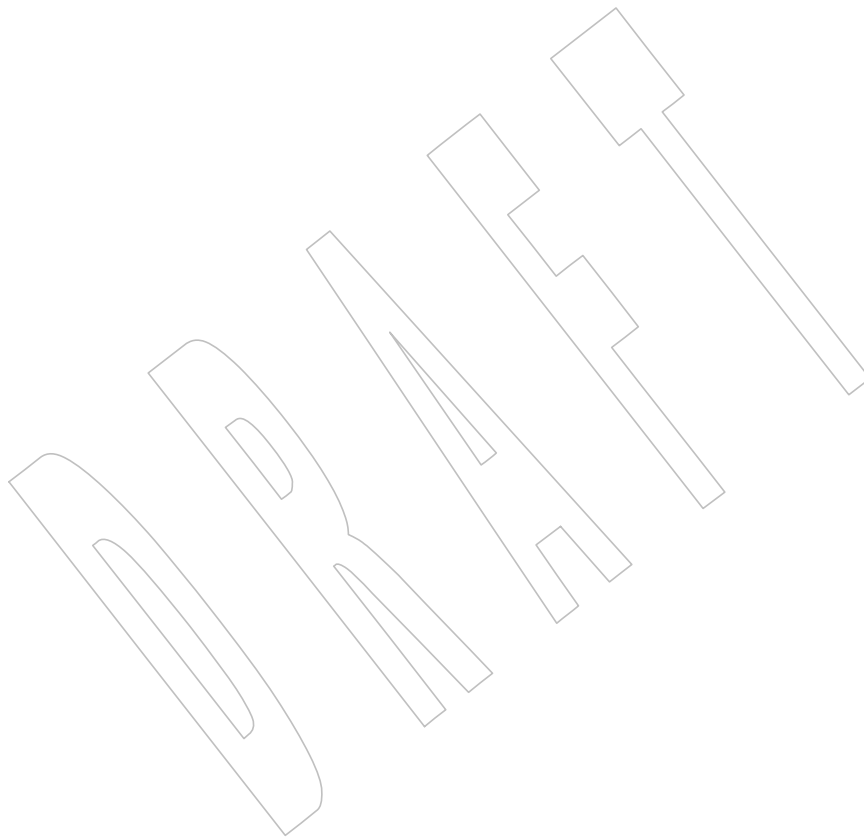
11882 `{MAX_CANON}`
 11883 For most purposes, the current minimum is adequate. Unless high-speed burst serial
 11884 devices are used, it should be left as is.

- 11885 {MAX_INPUT}
11886 See {MAX_CANON}.
- 11887 {MQ_OPEN_MAX}
11888 The current minimum should be adequate for most profiles.
- 11889 {MQ_PRIO_MAX}
11890 The current minimum corresponds to the required number of process scheduling priorities.
11891 Many realtime practitioners believe that the number of message priority levels ought to be
11892 the same as the number of execution scheduling priorities.
- 11893 {NAME_MAX}
11894 Many implementations now support larger values, and many applications and users
11895 assume that larger names can be used. Many existing profiles also specify a larger value.
11896 Specifying this value will reduce the number of conforming implementations, although this
11897 might not be a significant consideration over time. Values greater than 255 should not be
11898 required.
- 11899 {NGROUPS_MAX}
11900 The value selected will typically be 8 or larger.
- 11901 {OPEN_MAX}
11902 The historically common value for this has been 20. Many implementations support larger
11903 values. If applications that use larger values are anticipated, an appropriate value should be
11904 specified.
- 11905 {PAGESIZE}
11906 This is not actually a limit, but an implementation parameter. No profile should impose a
11907 requirement on this value.
- 11908 {PATH_MAX}
11909 Historically, the minimum has been either 1024 or indefinite, depending on the
11910 implementation. Few applications actually require values larger than 256, but some users
11911 may create file hierarchies that must be accessed with longer paths. This value should only
11912 be changed if there is a clear requirement.
- 11913 {PIPE_BUF}
11914 The current minimum is adequate for most applications. Historically, it has been larger. If
11915 applications that write single transactions larger than this are anticipated, it should be
11916 increased. Applications that write lines of text larger than this probably do not need it
11917 increased, as the text line is delimited by a <newline>.
- 11918 {POSIX_VERSION}
11919 This is actually not a limit, but a standard version stamp. Generally, a profile should specify
11920 IEEE Std 1003.1-200x by a name in the normative references section, not this value.
- 11921 {PTHREAD_DESTRUCTOR_ITERATIONS}
11922 It is unlikely that applications will need larger values to avoid loss of memory resources.
- 11923 {PTHREAD_KEYS_MAX}
11924 The current value should be adequate for most profiles.
- 11925 {PTHREAD_STACK_MIN}
11926 This should not be treated as an actual limit, but as an implementation parameter. No
11927 profile should impose a requirement on this value.
- 11928 {PTHREAD_THREADS_MAX}
11929 It is believed that most implementations will provide larger values.

- 11930 {RTSIG_MAX}
 11931 The current limit was chosen so that the set of POSIX.1 signal numbers can fit within a
 11932 32-bit field. It is recognized that most existing implementations define many more signals
 11933 than are specified in POSIX.1 and, in fact, many implementations have already exceeded 32
 11934 signals (including the “null signal”). Support of {_POSIX_RTSIG_MAX} additional signals
 11935 may push some implementations over the single 32-bit word line, but is unlikely to push
 11936 any implementations that are already over that line beyond the 64 signal line.
- 11937 {SEM_NSEMS_MAX}
 11938 The current value should be adequate for most profiles.
- 11939 {SEM_VALUE_MAX}
 11940 The current value should be adequate for most profiles.
- 11941 {SSIZE_MAX}
 11942 This limit reflects fundamental hardware characteristics (the size of an integer), and should
 11943 not be specified unless it is clearly required. Extreme care should be taken to assure that
 11944 any value that might be specified does not unnecessarily eliminate implementations
 11945 because of accidents of hardware design.
- 11946 {STREAM_MAX}
 11947 This limit is very closely related to {OPEN_MAX}. It should never be larger than
 11948 {OPEN_MAX}, but could reasonably be smaller for application areas where most files are
 11949 not accessed through *stdio*. Some implementations may limit {STREAM_MAX} to 20 but
 11950 allow {OPEN_MAX} to be considerably larger. Such implementations should be allowed for
 11951 if the applications permit.
- 11952 {TIMER_MAX}
 11953 The current limit should be adequate for most profiles, but it may need to be larger for
 11954 applications with a large number of asynchronous operations.
- 11955 {TTY_NAME_MAX}
 11956 This is not actually a limit, but an implementation parameter. No profile should impose a
 11957 requirement on this value.
- 11958 {TZNAME_MAX}
 11959 The minimum has been historically adequate, but if longer timezone names are anticipated
 11960 (particularly such values as UTC-1), this should be increased.

11961 D.3.6 Optional Behavior

11962 In IEEE Std 1003.1-200x, there are no instances of the terms unspecified, undefined,
 11963 implementation-defined, or with the verbs “may” or “need not”, that the developers of
 11964 IEEE Std 1003.1-200x anticipate or sanction as suitable for profile or test method citation. All of
 11965 these are merely warnings to conforming applications to avoid certain areas that can vary from
 11966 system to system, and even over time on the same system. In many cases, these terms are used
 11967 explicitly to support extensions, but profiles should not anticipate and require such extensions;
 11968 future versions of IEEE Std 1003.1 may do so.



11969

Rationale (Informative)

11970

Part E:

11971

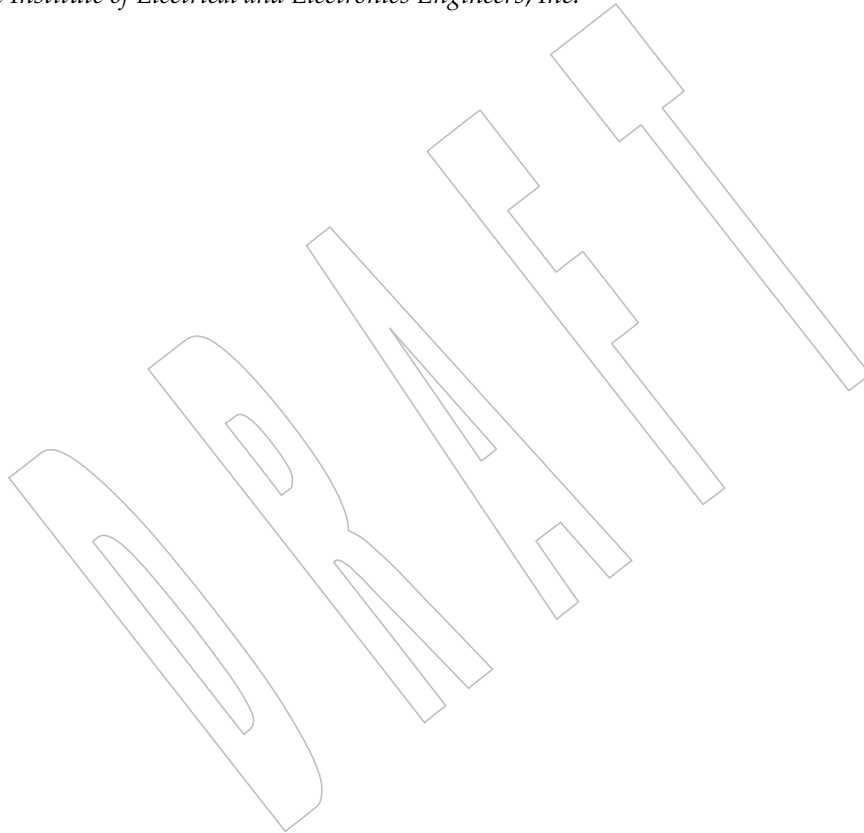
Subprofiling Considerations

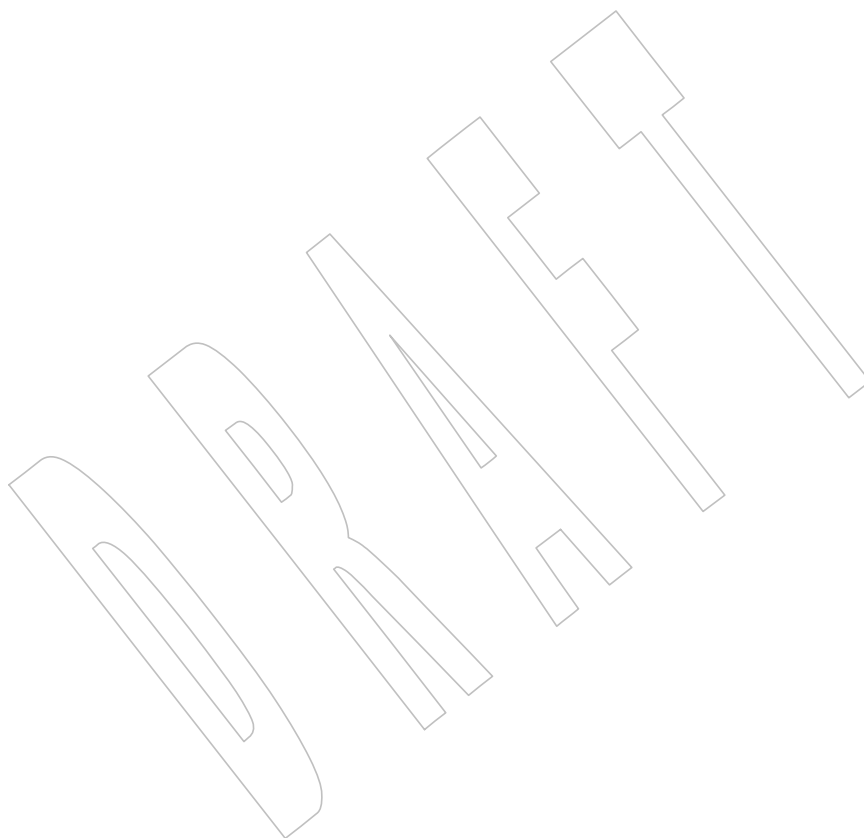
11972

The Open Group

11973

The Institute of Electrical and Electronics Engineers, Inc.





Subprofiling Considerations (Informative)

11974

11975

11976

11977

11978

11979

This section contains further information to satisfy the requirement that the project scope enable subprofiling of IEEE Std 1003.1-200x. The approach taken is to include a general requirement in normative text regarding subprofiling and to include an informative section (here) containing a proposed set of subprofiling options.

11980

E.1 Subprofiling Option Groups

11981

11982

11983

11984

11985

The following Option Groups⁴ are defined to support profiling. Systems claiming support to IEEE Std 1003.1-200x need not implement these options apart from the requirements stated in the Base Definitions volume of IEEE Std 1003.1-200x, Section 2.1.3, POSIX Conformance. These Option Groups allow profiles to subset the System Interfaces volume of IEEE Std 1003.1-200x by collecting sets of related functions.

11986

POSIX_ASYNCHRONOUS_IO: Asynchronous Input and Output Functions

11987

aio_cancel(), *aio_error()*, *aio_fsync()*, *aio_read()*, *aio_return()*, *aio_suspend()*, *aio_write()*,
lio_listio()

11988

11989

POSIX_BARRIERS: Barriers

11990

pthread_barrier_destroy(), *pthread_barrier_init()*, *pthread_barrier_wait()*, *pthread_barrierattr()*

11991

POSIX_C_LANG_JUMP: Jump Functions

11992

longjmp(), *setjmp()*

11993

POSIX_C_LANG_MATH: Maths Library

11994

acos(), *acosf()*, *acosh()*, *acoshf()*, *acoshl()*, *acosl()*, *asin()*, *asinf()*, *asinh()*, *asinhf()*, *asinhl()*,

11995

asinl(), *atan()*, *atan2()*, *atan2f()*, *atan2l()*, *atanf()*, *atanh()*, *atanhf()*, *atanhl()*, *atanl()*, *cabs()*,

11996

cabsf(), *cabsl()*, *cacos()*, *cacosf()*, *cacosh()*, *cacoshf()*, *cacoshl()*, *cacosl()*, *carg()*, *cargf()*, *cargl()*,

11997

casin(), *casinf()*, *casinh()*, *casinhf()*, *casinhl()*, *casinl()*, *catan()*, *catanf()*, *catanh()*, *catanhf()*,

11998

catanhl(), *catanl()*, *cbrt()*, *cbrtf()*, *cbrtl()*, *ccos()*, *ccosf()*, *ccosh()*, *ccoshf()*, *ccoshl()*, *ccosl()*,

11999

ceil(), *ceilf()*, *ceilL()*, *cexp()*, *cexpf()*, *cexpl()*, *cimag()*, *cimagf()*, *cimagl()*, *clog()*, *clogf()*, *clogl()*,

12000

conj(), *conjf()*, *conjl()*, *copysign()*, *copysignf()*, *copysignl()*, *cos()*, *cosf()*, *cosh()*, *coshf()*,

12001

coshl(), *cosl()*, *cpow()*, *cpowf()*, *cpowl()*, *cproj()*, *cprojf()*, *cprojl()*, *creal()*, *crealf()*, *creall()*,

12002

csin(), *csinf()*, *csinh()*, *csinhf()*, *csinhl()*, *csinl()*, *csqrt()*, *csqrtf()*, *csqrtl()*, *ctan()*, *ctanf()*,

12003

ctanh(), *ctanhf()*, *ctanhl()*, *ctanl()*, *erf()*, *erfc()*, *erfcf()*, *erfcl()*, *erff()*, *erfl()*, *exp()*, *exp2()*,

12004

exp2f(), *exp2l()*, *expf()*, *expl()*, *expm1()*, *expm1f()*, *expm1l()*, *fabs()*, *fabsf()*, *fabsl()*, *fdim()*,

12005

fdimf(), *fdiml()*, *floor()*, *floorf()*, *floorl()*, *fma()*, *fmaf()*, *fmal()*, *fmax()*, *fmaxf()*, *fmaxl()*, *fmin()*,

12006

fminf(), *fminl()*, *fmod()*, *fmodf()*, *fmodl()*, *fpclassify()*, *frexp()*, *frexpf()*, *frexpl()*, *hypot()*,

12007

hypotf(), *hypotl()*, *ilogb()*, *ilogbf()*, *ilogbl()*, *isfinite()*, *isgreater()*, *isgreaterequal()*, *isinf()*,

12008

isless(), *islessequal()*, *islessgreater()*, *isnan()*, *isnormal()*, *isunordered()*, *ldexp()*, *ldexpf()*,

12009

ldexpl(), *lgamma()*, *lgammaf()*, *lgammal()*, *llrint()*, *llrintf()*, *llrintl()*, *llround()*, *llroundf()*,

12010

llroundl(), *log()*, *log10()*, *log10f()*, *log10l()*, *log1p()*, *log1pf()*, *log1pl()*, *log2()*, *log2f()*, *log2l()*,

12011

logb(), *logbf()*, *logbl()*, *logf()*, *logl()*, *lrint()*, *lrintf()*, *lrintl()*, *lround()*, *lroundf()*, *lroundl()*,

12012

modf(), *modff()*, *modfl()*, *nan()*, *nanf()*, *nanl()*, *nearbyint()*, *nearbyintf()*, *nearbyintl()*,

12013

nextafter(), *nextafterf()*, *nextafterl()*, *nexttoward()*, *nexttowardf()*, *nexttowardl()*, *pow()*, *powf()*,

12014

4. These are modelled on the Units of Functionality from IEEE Std 1003.13-1998.

12015 *powl()*, *remainder()*, *remainderf()*, *remainderl()*, *remquo()*, *remquof()*, *remquol()*, *rint()*, *rintf()*,
 12016 *rintl()*, *round()*, *roundf()*, *roundl()*, *scalbln()*, *scalblnf()*, *scalblnl()*, *scalbn()*, *scalbnf()*,
 12017 *scalbnl()*, *signbit()*, *sin()*, *sinf()*, *sinh()*, *sinhf()*, *sinhl()*, *sinl()*, *sqrt()*, *sqrtf()*, *sqrtl()*, *tan()*,
 12018 *tanf()*, *tanh()*, *tanhf()*, *tanhL()*, *tanl()*, *tgamma()*, *tgammaf()*, *tgammaL()*, *trunc()*, *truncf()*,
 12019 *truncl()*

12020 POSIX_C_LANG_SUPPORT: General ISO C Library

12021 *abs()*, *asctime()*, *atof()*, *atoi()*, *atol()*, *atoll()*, *bsearch()*, *calloc()*, *ctime()*, *difftime()*, *div()*,
 12022 *feclearexcept()*, *fegetenv()*, *fegetexceptflag()*, *fegetround()*, *fehldexcept()*, *feraiseexcept()*,
 12023 *fesetenv()*, *fesetexceptflag()*, *fesetround()*, *fetestexcept()*, *feupdateenv()*, *free()*, *gmtime()*,
 12024 *imaxabs()*, *imaxdiv()*, *isalnum()*, *isalpha()*, *isblank()*, *iscntrl()*, *isdigit()*, *isgraph()*, *islower()*,
 12025 *isprint()*, *ispunct()*, *isspace()*, *isupper()*, *isxdigit()*, *labs()*, *ldiv()*, *llabs()*, *lldiv()*, *localeconv()*,
 12026 *localtime()*, *malloc()*, *memchr()*, *memcmp()*, *memcpy()*, *memmove()*, *memset()*, *mktime()*,
 12027 *qsort()*, *rand()*, *realloc()*, *setlocale()*, *snprintf()*, *sprintf()*, *srand()*, *sscanf()*, *strcat()*, *strchr()*,
 12028 *strcmp()*, *strcoll()*, *strcpy()*, *strcspn()*, *strerror()*, *strftime()*, *strlen()*, *strncat()*, *strncpy()*,
 12029 *strncpy()*, *strpbrk()*, *strrchr()*, *strspn()*, *strstr()*, *strtod()*, *strtof()*, *strtoimax()*, *strtok()*, *strtol()*,
 12030 *strtol()*, *strtoll()*, *strtoul()*, *strtoull()*, *strtoumax()*, *strxfrm()*, *time()*, *tolower()*, *toupper()*,
 12031 *tzname*, *tzset()*, *va_arg()*, *va_copy()*, *va_end()*, *va_start()*, *vsprintf()*, *vsscanf()*

12032 POSIX_C_LANG_SUPPORT_R: Thread-Safe General ISO C Library

12033 *asctime_r()*, *ctime_r()*, *gmtime_r()*, *localtime_r()*, *rand_r()*, *strerror_r()*, *strtok_r()*

12034 POSIX_C_LANG_WIDE_CHAR: Wide-Character ISO C Library

12035 *btowc()*, *iswalnum()*, *iswalphalower()*, *iswblank()*, *iswcntrl()*, *iswctype()*, *iswdigit()*, *iswgraph()*,
 12036 *iswlower()*, *iswprint()*, *iswpunct()*, *iswspace()*, *iswupper()*, *iswxdigit()*, *mblen()*, *mbrlen()*,
 12037 *mbrtowc()*, *mbsinit()*, *mbsrtowcs()*, *mbstowcs()*, *mbtowc()*, *swprintf()*, *swscanf()*, *towctrans()*,
 12038 *towlower()*, *towupper()*, *vwprintf()*, *vwscanf()*, *wcrtomb()*, *wcscat()*, *wcschr()*, *wcscmp()*,
 12039 *wcscoll()*, *wcscpy()*, *wcscspn()*, *wcsftime()*, *wcslen()*, *wcsncat()*, *wcsncpy()*, *wcsncpy()*,
 12040 *wcspbrk()*, *wcsrchr()*, *wcstombs()*, *wcsspn()*, *wcsstr()*, *wcstod()*, *wcstof()*, *wcstoimax()*,
 12041 *wcstok()*, *wcstol()*, *wcstold()*, *wcstoll()*, *wcstombs()*, *wcstoul()*, *wcstoull()*, *wcstoumax()*,
 12042 *wcsxfrm()*, *wctob()*, *wctomb()*, *wctrans()*, *wctype()*, *wmemchr()*, *wmemcmp()*, *wmemcpy()*,
 12043 *wmemmove()*, *wmemset()*

12044 POSIX_C_LANG_WIDE_CHAR_EXT: Extended Wide-Character ISO C Library

12045 *mbsnrtowcs()*, *wcpcpy()*, *wcpncpy()*, *wcscasecmp()*, *wcsdup()*, *wcsncasecmp()*, *wcsnlen()*,
 12046 *wcsnrtombs()*

12047 POSIX_C_LIB_EXT: General C Library Extension

12048 *fnmatch()*, *getopt()*, *getsubopt()*, *optarg*, *opterr*, *optind*, *optopt*, *stpcpy()*, *stpncpy()*, *strcasecmp()*,
 12049 *strdup()*, *strfmon()*, *strncasecmp()*, *strndup()*, *strlen()*

12050 POSIX_CLOCK_SELECTION: Clock Selection

12051 *clock_nanosleep()*, *pthread_condattr_getclock()*, *pthread_condattr_setclock()*

12052 POSIX_DEVICE_IO: Device Input and Output

12053 *FD_CLR()*, *FD_ISSET()*, *FD_SET()*, *FD_ZERO()*, *clearerr()*, *close()*, *fclose()*, *fdopen()*, *feof()*,
 12054 *ferror()*, *fflush()*, *fgetc()*, *fgets()*, *fileno()*, *fopen()*, *fprintf()*, *fputc()*, *fputs()*, *fread()*, *freopen()*,
 12055 *fscanf()*, *fwrite()*, *getc()*, *getchar()*, *gets()*, *open()*, *perror()*, *poll()*, *printf()*, *pread()*, *pselect()*,
 12056 *putc()*, *putchar()*, *puts()*, *pwrite()*, *read()*, *scanf()*, *select()*, *setbuf()*, *setvbuf()*, *stderr*, *stdin*,
 12057 *stdout*, *ungetc()*, *vfprintf()*, *vscanf()*, *vprintf()*, *vscanf()*, *write()*

12058 POSIX_DEVICE_IO_EXT: Extended Device Input and Output

12059 *dprintf()*, *fmemopen()*, *open_memstream()*

12060 POSIX_DEVICE_SPECIFIC: General Terminal

12061 *cfgetispeed()*, *cfgetospeed()*, *cfsetispeed()*, *cfsetospeed()*, *ctermid()*, *isatty()*, *tcdrain()*, *tcflow()*,
 12062 *tcflush()*, *tcgetattr()*, *tcsendbreak()*, *tcsetattr()*, *ttyname()*

12063	POSIX_DEVICE_SPECIFIC_R: Thread-Safe General Terminal
12064	<i>ttyname_r()</i>
12065	POSIX_DYNAMIC_LINKING: Dynamic Linking
12066	<i>dlclose(), dlerror(), dlopen(), dlsym()</i>
12067	POSIX_FD_MGMT: File Descriptor Management
12068	<i>dup(), dup2(), fcntl(), fgetpos(), fseek(), fseeko(), fsetpos(), ftell(), ftello(), ftruncate(), lseek(),</i>
12069	<i>rewind()</i>
12070	POSIX_FIFO: FIFO
12071	<i>mkfifo()</i>
12072	POSIX_FIFO_FD: FIFO File Descriptor Routines
12073	<i>mkfifoat(), mknodat()</i>
12074	POSIX_FILE_ATTRIBUTES: File Attributes
12075	<i>chmod(), chown(), fchmod(), fchown(), umask()</i>
12076	POSIX_FILE_ATTRIBUTES_FD: File Attributes File Descriptor Routines
12077	<i>fchmodat(), fchownat()</i>
12078	POSIX_FILE_LOCKING: Thread-Safe Stdio Locking
12079	<i>flockfile(), frylockfile(), funlockfile(), getc_unlocked(), getchar_unlocked(), putc_unlocked(),</i>
12080	<i>putchar_unlocked()</i>
12081	POSIX_FILE_SYSTEM: File System
12082	<i>access(), chdir(), closedir(), creat(), fchdir(), fpathconf(), fstat(), fstatvfs(), getcwd(), link(),</i>
12083	<i>mkdir(), mkstemp(), opendir(), pathconf(), readdir(), remove(), rename(), rewinddir(), rmdir(),</i>
12084	<i>stat(), statvfs(), tmpfile(), tmpnam(), truncate(), unlink(), utime()</i>
12085	POSIX_FILE_SYSTEM_EXT: File System Extensions
12086	<i>alphasort(), dirfd(), getdelim(), getline(), mkdtemp(), scandir()</i>
12087	POSIX_FILE_SYSTEM_FD: File System File Descriptor Routines
12088	<i>faccessat(), fdopendir(), fstatat(), futimesat(), linkat(), mkdirat(), openat(), renameat(),</i>
12089	<i>unlinkat()</i>
12090	POSIX_FILE_SYSTEM_GLOB: File System Glob Expansion
12091	<i>glob(), globfree()</i>
12092	POSIX_FILE_SYSTEM_R: Thread-Safe File System
12093	<i>readdir_r()</i>
12094	POSIX_I18N: Internationalization
12095	<i>catclose(), catgets(), catopen(), iconv(), iconv_close(), iconv_open(), nl_langinfo()</i>
12096	POSIX_JOB_CONTROL: Job Control
12097	<i>setpgid(), tcgetpgrp(), tcsetpgrp(), tcgetsid()</i>
12098	POSIX_MAPPED_FILES: Memory Mapped Files
12099	<i>mmap(), munmap()</i>
12100	POSIX_MEMORY_PROTECTION: Memory Protection
12101	<i>mprotect()</i>
12102	POSIX_MULTI_CONCURRENT_LOCALES: Multiple Concurrent Locales
12103	<i>duplocale(), freelocale(), isalnum_l(), isalpha_l(), isblank_l(), iscntrl_l(), isdigit_l(), isgraph_l(),</i>
12104	<i>islower_l(), isprint_l(), ispunct_l(), isspace_l(), isupper_l(), iswalnum_l(), iswalalpha_l(),</i>
12105	<i>iswblank_l(), iswcntrl_l(), iswctype_l(), iswdigit_l(), iswgraph_l(), iswlower_l(), iswprint_l(),</i>
12106	<i>iswpunct_l(), iswspace_l(), iswupper_l(), iswxdigit_l(), isxdigit_l(), newlocale(), strcasecmp_l(),</i>

12107	<i>strcoll_1()</i> , <i>strfmon_1()</i> , <i>strncasecmp_1()</i> , <i>strxfrm_1()</i> , <i>tolower_1()</i> , <i>toupper_1()</i> , <i>towctrans_1()</i> ,
12108	<i>towlower()</i> , <i>towupper()</i> , <i>uselocale()</i> , <i>wcscasecmp_1()</i> , <i>wcscoll_1()</i> , <i>wcsncasecmp_1()</i> , <i>wcsxfrm_1()</i> ,
12109	<i>wctrans_1()</i> , <i>wctype_1()</i>
12110	POSIX_MULTI_PROCESS: Multiple Processes
12111	<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> ,
12112	<i>exit()</i> , <i>fork()</i> , <i>getpgrp()</i> , <i>getpgid()</i> , <i>getpid()</i> , <i>getppid()</i> , <i>getsid()</i> , <i>setsid()</i> , <i>sleep()</i> , <i>times()</i> , <i>wait()</i> ,
12113	<i>waitid()</i> , <i>waitpid()</i>
12114	POSIX_MULTI_PROCESS_FD: Multiple Processes File Descriptor Routines
12115	<i>fexecve()</i>
12116	POSIX_NETWORKING: Networking
12117	<i>accept()</i> , <i>bind()</i> , <i>connect()</i> , <i>endhostent()</i> , <i>endnetent()</i> , <i>endprotoent()</i> , <i>endservent()</i> ,
12118	<i>freeaddrinfo()</i> , <i>gai_strerror()</i> , <i>getaddrinfo()</i> , <i>gethostent()</i> , <i>gethostname()</i> , <i>getnameinfo()</i> ,
12119	<i>getnetbyaddr()</i> , <i>getnetbyname()</i> , <i>getnetent()</i> , <i>getpeername()</i> , <i>getprotobyname()</i> ,
12120	<i>getprotobynumber()</i> , <i>getprotoent()</i> , <i>getserbyname()</i> , <i>getservbyport()</i> , <i>getservent()</i> ,
12121	<i>getsockname()</i> , <i>getsockopt()</i> , <i>htonl()</i> , <i>htons()</i> , <i>if_freenameindex()</i> , <i>if_indextoname()</i> ,
12122	<i>if_nameindex()</i> , <i>if_nametoindex()</i> , <i>inet_addr()</i> , <i>inet_ntoa()</i> , <i>inet_ntop()</i> , <i>inet_pton()</i> , <i>listen()</i> ,
12123	<i>ntohl()</i> , <i>ntohs()</i> , <i>recv()</i> , <i>recvfrom()</i> , <i>recvmsg()</i> , <i>send()</i> , <i>sendmsg()</i> , <i>sendto()</i> , <i>sethostent()</i> ,
12124	<i>setnetent()</i> , <i>setprotoent()</i> , <i>setservent()</i> , <i>setsockopt()</i> , <i>shutdown()</i> , <i>socket()</i> , <i>socketmark()</i> ,
12125	<i>socketpair()</i>
12126	POSIX_PIPE: Pipe
12127	<i>pipe()</i>
12128	POSIX_ROBUST_MUTEXES: Robust Mutexes
12129	<i>pthread_mutex_consistent()</i> , <i>pthread_mutexattr_getrobust()</i> , <i>pthread_mutexattr_setrobust()</i>
12130	POSIX_REALTIME_SIGNALS: Realtime Signals
12131	<i>sigqueue()</i> , <i>sigtimedwait()</i> , <i>sigwaitinfo()</i>
12132	POSIX_REGEX: Regular Expressions
12133	<i>regcomp()</i> , <i>regerror()</i> , <i>regexec()</i> , <i>regfree()</i>
12134	POSIX_RW_LOCKS: Reader Writer Locks
12135	<i>pthread_rwlock_destroy()</i> , <i>pthread_rwlock_init()</i> , <i>pthread_rwlock_rdlock()</i> ,
12136	<i>pthread_rwlock_timedrdlock()</i> , <i>pthread_rwlock_timedwrlock()</i> , <i>pthread_rwlock_tryrdlock()</i> ,
12137	<i>pthread_rwlock_trywrlock()</i> , <i>pthread_rwlock_unlock()</i> , <i>pthread_rwlock_wrlock()</i> ,
12138	<i>pthread_rwlockattr_destroy()</i> , <i>pthread_rwlockattr_init()</i> , <i>pthread_rwlockattr_getpshared()</i> ,
12139	<i>pthread_rwlockattr_setpshared()</i>
12140	POSIX_SEMAPHORES: Semaphores
12141	<i>sem_close()</i> , <i>sem_destroy()</i> , <i>sem_getvalue()</i> , <i>sem_init()</i> , <i>sem_open()</i> , <i>sem_post()</i> ,
12142	<i>sem_timedwait()</i> , <i>sem_trywait()</i> , <i>sem_unlink()</i> , <i>sem_wait()</i>
12143	POSIX_SHELL_FUNC: Shell and Utilities
12144	<i>pclose()</i> , <i>popen()</i> , <i>system()</i> , <i>wordexp()</i> , <i>wordfree()</i>
12145	POSIX_SIGNAL_JUMP: Signal Jump Functions
12146	<i>siglongjmp()</i> , <i>sigsetjmp()</i>
12147	POSIX_SIGNALS: Signals
12148	<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> ,
12149	<i>sigfillset()</i> , <i>sigismember()</i> , <i>signal()</i> , <i>sigpending()</i> , <i>sigprocmask()</i> , <i>sigsuspend()</i> , <i>sigwait()</i>
12150	POSIX_SIGNALS_EXT: Extended Signals
12151	<i>psignal()</i> , <i>psiginfo()</i> , <i>strsignal()</i>

POSIX_SINGLE_PROCESS: Single Process

confstr(), *environ*, *errno*, *getenv()*, *setenv()*, *sysconf()*, *uname()*, *unsetenv()*

POSIX_SPIN_LOCKS: Spin Locks

pthread_spin_destroy(), *pthread_spin_init()*, *pthread_spin_lock()*, *pthread_spin_trylock()*,
pthread_spin_unlock()

POSIX_SYMBOLIC_LINKS: Symbolic Links

lchown(),⁵ *lstat()*, *readlink()*, *symlink()*

POSIX_SYMBOLIC_LINKS_FD: Symbolic Links File Descriptor Routines

readlinkat(), *symlinkat()*

POSIX_SYSTEM_DATABASE: System Database

getgrgid(), *getgrnam()*, *getpwnam()*, *getpwuid()*

POSIX_SYSTEM_DATABASE_R: Thread-Safe System Database

getgrgid_r(), *getgrnam_r()*, *getpwnam_r()*, *getpwuid_r()*

POSIX_THREADS_BASE: Base Threads

pthread_atfork(), *pthread_attr_destroy()*, *pthread_attr_getdetachstate()*,
pthread_attr_getschedparam(), *pthread_attr_init()*, *pthread_attr_setdetachstate()*,
pthread_attr_setschedparam(), *pthread_cancel()*, *pthread_cleanup_pop()*, *pthread_cleanup_push()*,
pthread_cond_broadcast(), *pthread_cond_destroy()*, *pthread_cond_init()*, *pthread_cond_signal()*,
pthread_cond_timedwait(), *pthread_cond_wait()*, *pthread_condattr_destroy()*,
pthread_condattr_init(), *pthread_create()*, *pthread_detach()*, *pthread_equal()*, *pthread_exit()*,

12198	<i>setstate()</i> , <i>signgam</i> , <i>srand48()</i> , <i>srandom()</i> , <i>strptime()</i> , <i>swab()</i> , <i>tdelete()</i> , <i>tfind()</i> , <i>timezone()</i> ,
12199	<i>toascii()</i> , <i>tsearch()</i> , <i>twalk()</i>
12200	XSI_DBM: XSI Database Management
12201	<i>dbm_clearerr()</i> , <i>dbm_close()</i> , <i>dbm_delete()</i> , <i>dbm_error()</i> , <i>dbm_fetch()</i> , <i>dbm_firstkey()</i> ,
12202	<i>dbm_nextkey()</i> , <i>dbm_open()</i> , <i>dbm_store()</i>
12203	XSI_DEVICE_IO: XSI Device Input and Output
12204	<i>fntmsg()</i> , <i>readv()</i> , <i>writv()</i>
12205	XSI_DEVICE_SPECIFIC: XSI General Terminal
12206	<i>grantpt()</i> , <i>posix_openpt()</i> , <i>ptsname()</i> , <i>unlockpt()</i>
12207	XSI_FILE_SYSTEM: XSI File System
12208	<i>basename()</i> , <i>dirname()</i> , <i>ftw()</i> , <i>lockf()</i> , <i>mknod()</i> , <i>nftw()</i> , <i>realpath()</i> , <i>seekdir()</i> , <i>sync()</i> , <i>telldir()</i> ,
12209	<i>tempnam()</i>
12210	XSI_IPC: XSI Interprocess Communication
12211	<i>ftok()</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> ,
12212	<i>shmdt()</i> , <i>shmget()</i>
12213	XSI_JUMP: XSI Jump Functions
12214	<i>_longjmp()</i> , <i>_setjmp()</i>
12215	XSI_MATH: XSI Maths Library
12216	<i>j0()</i> , <i>j1()</i> , <i>jn()</i> , <i>y0()</i> , <i>y1()</i> , <i>yn()</i>
12217	XSI_MULTI_PROCESS: XSI Multiple Process
12218	<i>getpriority()</i> , <i>getrlimit()</i> , <i>getrusage()</i> , <i>nice()</i> , <i>setpgrp()</i> , <i>setpriority()</i> , <i>setrlimit()</i> , <i>ulimit()</i> ,
12219	XSI_SIGNALS: XSI Signal
12220	<i>killpg()</i> , <i>sigaltstack()</i> , <i>sighold()</i> , <i>sigignore()</i> , <i>siginterrupt()</i> , <i>sigpause()</i> , <i>sigrelse()</i> , <i>sigset()</i> ,
12221	XSI_SINGLE_PROCESS: XSI Single Process
12222	<i>gethostid()</i> , <i>gettimeofday()</i> , <i>putenv()</i>
12223	XSI_SYSTEM_DATABASE: XSI System Database
12224	<i>endpwent()</i> , <i>getpwent()</i> , <i>setpwent()</i>
12225	XSI_SYSTEM_LOGGING: XSI System Logging
12226	<i>closelog()</i> , <i>openlog()</i> , <i>setlogmask()</i> , <i>syslog()</i>
12227	XSI_THREADS_EXT: XSI Threads Extensions
12228	<i>pthread_attr_getstack()</i> , <i>pthread_attr_setstack()</i> , <i>pthread_getconcurrency()</i> ,
12229	<i>pthread_setconcurrency()</i>
12230	XSI_TIMERS: XSI Timers
12231	<i>getitimer()</i> , <i>setitimer()</i>
12232	XSI_USER_GROUPS: XSI User and Group
12233	<i>endgrent()</i> , <i>endutxent()</i> , <i>getgrent()</i> , <i>getutxent()</i> , <i>getutxid()</i> , <i>getutxline()</i> , <i>pututxline()</i> ,
12234	<i>setgrent()</i> , <i>setregid()</i> , <i>setreuid()</i> , <i>setutxent()</i>
12235	XSI_WIDE_CHAR: XSI Wide-Character Library
12236	<i>wcswidth()</i> , <i>wcwidth()</i>

Index

/dev/tty.....	21	_POSIX_THREAD_PRIORITY_SCHEDULING.....	284
/etc/passwd	34	_POSIX_THREAD_PRIO_INHERIT.....	284
<pthread.h>.....	160	_POSIX_THREAD_PRIO_PROTECT	284
_asm_builtin_atoi().....	84	_POSIX_THREAD_PROCESS_SHARED.....	284
_exit()	100, 118	_POSIX_THREAD_SAFE_FUNCTIONS	10, 284
_Exit().....	271	_POSIX_THREAD_SPORADIC_SERVER	284
_exit().....	271	_POSIX_TIMEOUTS	10, 284
_longjmp().....	274	_POSIX_TIMERS	10, 284
_POSIX_ADVISORY_INFO	280	_POSIX_TRACE.....	284
_POSIX_ASYNCHRONOUS_IO.....	10, 280	_POSIX_TRACE_EVENT_FILTER.....	284
_POSIX_BARRIERS.....	10, 280	_POSIX_TRACE_INHERIT.....	284
_POSIX_CHOWN_RESTRICTED.....	5, 281	_POSIX_TRACE_LOG	285
_POSIX_CLOCK_SELECTION	10, 281	_POSIX_TYPED_MEMORY_OBJECTS.....	285
_POSIX_CPUTIME.....	281	_POSIX_TZNAME_MAX.....	58
_POSIX_C_SOURCE	85, 88	_POSIX_VDISABLE	5
_POSIX_FSYNC	281	SC_PAGESIZE	117, 119
_POSIX_IPV6	281	setjmp().....	274
_POSIX_JOB_CONTROL	5, 281, 285	_XOPEN_SOURCE.....	85
_POSIX_MAPPED_FILES	10, 281	__errno().....	92
_POSIX_MEMLOCK.....	281	abort().....	271
_POSIX_MEMLOCK_RANGE	282	access	272
_POSIX_MEMORY_PROTECTION.....	10, 282	access().....	34
_POSIX_MESSAGE_PASSING	282	active trace stream.....	202
_POSIX_MONOTONIC_CLOCK.....	282	adb, rationale for omission	260
_POSIX_NO_TRUNC.....	5	address families	177
_POSIX_PRIORITIZED_IO.....	282	addressing	177
_POSIX_PRIORITY_SCHEDULING	282	advisory information	104
_POSIX_READER_WRITER_LOCKS.....	10	aio_cancel().....	112-113
_POSIX_REALTIME_SIGNALS.....	10, 282	aio_fsync().....	97, 111
_POSIX_REGEX	282	AIO_LISTIO_MAX.....	285
_POSIX_SIG_MAX.....	94, 287	AIO_MAX.....	285
_POSIX_SAVED_IDS	5, 282	AIO_PRIO_DELTA_MAX	285
_POSIX_SEMAPHORES.....	10, 283	aio_read()	113
_POSIX_SHARED_MEMORY_OBJECTS	283	aio_suspend().....	111, 138
_POSIX_SHELL	283	aio_write().....	113
_POSIX_SOURCE.....	85	alarm.....	271
_POSIX_SPAWN.....	283	alarm()	102, 137
_POSIX_SPINLOCKS.....	283	alias.....	229, 275
_POSIX_SPIN_LOCKS.....	10	alias substitution.....	232
_POSIX_SPORADIC_SERVER	283	AND lists	248
_POSIX_SS_REPL_MAX	133	application instrumentation	190
_POSIX_SYNCHRONIZED_IO.....	283	application-managed thread stack.....	176
_POSIX_SYNC_IO.....	281	appropriate privilege	13, 25
_POSIX_THREADS.....	10, 283	ar.....	277-278
_POSIX_THREAD_ATTR_STACKADDR.....	283	arbitrary file size.....	228
_POSIX_THREAD_ATTR_STACKSIZE.....	283	ARG_MAX	23, 222, 285
		arithmetic expansion.....	238

- as, rationale for omission260
- asa275, 277, 279
- ASCII24
- async-cancel safety174
- async-signal-safe100
- asynchronous error178
- asynchronous I/O111, 272
- asynchronous lists247
- at275
- atexit()174
- atoi()83-84
- awk275, 277
- background19-21, 68
- backslash230
- banner, rationale for omission260
- barriers153
- basename274, 276
- basic regular expression61
- batch275
 - general concepts257
- batch environment254
 - option definitions255
- batch environment utilities
 - common behavior260
- batch services259
- batch systems
 - historical implementations254
 - history254
- bc274-275, 280
- bcmp205
- bcopy205
- BC_BASE_MAX222
- BC_DIM_MAX222
- BC_SCALE_MAX222
- bg229, 275
- bounded response273
- bracket expression
 - grammar65
- BRE
 - expression anchoring63
 - grammar lexical conventions65
 - matching a collating element61
 - matching a single character61
 - matching multiple characters62
 - ordinary character61
 - periods61
 - precedence63
 - special character61
- BSD16, 18, 21, 24, 67-69, 93, 95, 99-100, 203
- bsd_signal205
- built-in utilities228
- bzero205
- C Shell18-20
- C-language extensions269, 274
- c99277
- calendar, rationale for omission260
- cancel, rationale for omission260
- cancellation cleanup handler172, 174
- cancellation cleanup stack172
- canonical mode input processing69
- case248
- case folding34-35
- cat228
- catclose()276
- catgets()276
- catopen()276
- cd229, 276
- CEO62
- change history79, 219
- char type204
- character14
- character encoding43
- character set42
- character set description file43
- character set, portable filename24
- character, rationale14
- CHARCLASS_NAME_MAX48
- CHAR_MAX50
- chgrp228, 276-277
- child process14
- CHILD_MAX5, 204, 222, 247, 285
- chmod228, 272, 276-277
- chmod()26
- chown228, 272, 276-277
- chown()26
- chroot()25
- chroot, rationale for omission260
- cksum228, 276
- clock133
- clock tick14
- clock tick, rationale14
- CLOCKRES_MIN285
- clocks133
- clock_getcpuclockid()140-141
- clock_nanosleep()138
- CLOCK_PROCESS_CPUTIME_ID139, 141
- CLOCK_REALTIME133-138
- CLOCK_THREAD_CPUTIME_ID139, 141
- close272
- close()118
- closedir272
- closelog()277
- cmp228, 275
- codes8, 83, 221
- col, rationale for omission260

Index

collating element order.....	62
COLL_WEIGHTS_MAX.....	222
column position.....	15
COLUMNS.....	57
comm.....	275
command.....	14, 229, 274
command execution.....	245
command language.....	269, 274
command search.....	245
command substitution.....	237
compilation environment.....	84
complex data manipulation.....	269, 275
compound commands.....	248
concurrent execution.....	33
conditional construct	
case.....	248
if.....	249
configurable limits.....	280, 285
configuration interrogation.....	268, 271
configuration options.....	278
shell and utilities.....	278
system interfaces.....	280
conformance.....	6, 9, 12-13, 36, 79, 202, 219
conformance document.....	6
conformance document, rationale.....	6
conforming application.....	12, 93, 225, 227
conforming application, strictly.....	9, 12, 100
confstr.....	272
connection indication queue.....	178
control mode.....	71
controlling terminal.....	15, 68, 271
core.....	34
covert channel.....	36
cp.....	228, 276
cpio, rationale for omission.....	260
cpp, rationale for omission.....	260
creat().....	118, 228
crontab.....	275
CSIZE.....	71
csplit.....	275
ctags.....	277
ctime().....	274
cu, rationale for omission.....	260
cut.....	275
data access.....	268, 272
data type.....	202
date.....	276
dc, rationale for omission.....	261
dd.....	228, 275-276
defined types.....	202
definitions.....	83, 220
DELAYTIMER_MAX.....	285
determinism.....	268
device number.....	16
device, logical.....	22
df.....	228, 276-277
diff.....	275-276
dircmp, rationale for omission.....	261
direct I/O.....	16
directory.....	16
directory device.....	65
directory entry.....	16
directory files.....	65
directory protection.....	33
directory structure.....	65
directory, root.....	25
dirname.....	274, 276
dis, rationale for omission.....	261
display.....	16
dot.....	16
dot-dot.....	16, 24, 39
double-quotes.....	230
du.....	228, 276-277
dup.....	272
dup().....	118
dup2.....	272
dup2().....	118
EBUSY.....	91, 161
ECANCELED.....	90
echo.....	274
ECHOE.....	71
ECHOK.....	71
ECHONL.....	71
ecvt.....	205
ed.....	275-276
EDOM.....	91
EFAULT.....	89-90
effective user ID.....	34
EFTYPE.....	90
EILSEQ.....	91
EINPROGRESS.....	112
EINTR.....	90, 92, 102
EINVAL.....	90
ELOOP.....	90
emacs, rationale for omission.....	261
ENAMETOOLONG.....	90
endgrent().....	32
endpwent().....	32
ENOMEM.....	91
ENOSYS.....	91, 115
ENOTSUP.....	91
ENOTTY.....	66, 90-91
env.....	274, 277
environment access.....	268, 273
environment variable.....	56

definition	56	feature test macro	84-85, 203
E_OVERFLOW	91	fg	229, 275
EPERM	170	fgetc	272
EPIPE	91	fgetpos()	203
Epoch	17, 40, 134	field splitting	241
ERANGE	91	FIFO	17, 24, 127
ERASE	69	FIFO special file	17
ERE	63	file	17
alternation	64	file access permissions	34
bracket expression	64	file classes	17
expression anchoring	64	file creation	220
grammar	65	file descriptors	102
grammar lexical conventions	65	file format notation	42
matching a collating element	64	file hierarchy	34
matching a single character	64	file hierarchy manipulation	269, 276
matching multiple characters	64	file permissions	34, 68
ordinary character	64	file read	220
periods	64	file removal	220
precedence	64	file size, arbitrary	228
special character	64	file system	17
EROFS	91	file system, mounted	22
errno	89	file system, root	25
per-thread	91	file times update	36
error conditions	41	file write	220
error handling	252	file, passwd	24
error numbers	89, 93	filename	17, 34
errors	243	filename portability	36
escape character	230	fileno()	23
ex	275-276	filtering of trace event types	199
exec	196	filtering trace event types	199
exec family	19, 116, 174, 223, 228, 250, 271	find	275-276
Execution Time Monitoring	139	FIPS requirements	5
execution time, measurement	36	flockfile()	92
exit status	243	fnmatch()	277
exit()	100, 271	fold	275
EXIT_FAILURE	273	fopen	272
EXIT_SUCCESS	273	fopen()	18, 228
expand	275	for loop	248
expr	274-275	foreground	19-21, 67-68
EXPR_NEST_MAX	222	fork()	19, 67, 107, 116, 118, 196, 203, 207, 271
extended regular expression	63	fort77	277
extended security controls	34	fpathconf	272
false	229, 274	fpathconf()	271
fc	229, 275	fputc	272
fchmod	272	fread	272
fclose	272	free()	101, 173
fcntl	272	fseek	272
fcntl()	66, 90, 118, 121, 203	fseeko()	203
fcntl() locks	176	fsetpos	272
fcvt	205	fsetpos()	203
fdopen()	118	fstat	271-272
FD_CLOEXEC	121	fsync()	111
		ftello()	203

Index

ftime	205
ftruncate	272
ftruncate()	117-118, 120
ftw()	228
function definition command	249
functions	
implementation of	83
use of	83
fwrite	272
gcvrt	205
gencat	276
general terminal interface	66
getc	272
getc()	164
getch	272
getconf	221, 271, 277
getcontext	206
getegid	271
geteuid	271
getgid	271
getgrent()	32
getgrgid	271
getgrgid()	32, 165
getgrnam	271
getgrnam()	32, 36, 165
getgroups()	26
gethostbyaddr	206
gethostbyname	206
getlogin	271
getopt()	73, 276-277
getopts	229, 277
getpgrp()	20
getpid	271
getpid()	102
getppid	271
getpriority()	128
getpwent()	32
getpwnam	271
getpwnam()	32, 36, 165
getpwuid	271
getpwuid()	32, 165
getrlimit()	228
getrusage()	141
getty	68
getuid	271
getuid()	102, 204
getwd	206
gid_t	32
glob()	277
globfree()	277
gmtime()	40
grammar conventions	224
grep	275-276
group database	17
group database access	33
group file	18
grouping commands	248
head	275
headers	75
here-document	243
historical implementations	18
HOME	5
host byte order	36
hosted implementation	18
h_errno	206
ICANON	69, 71
iconv()	276
iconv_close()	276
iconv_open()	276
id	277
if	249
implementation	18
implementation, historical	18
implementation, hosted	18
implementation, native	23
implementation, specific	18
implementation-defined	6-7
implementation-defined, rationale	6
incomplete pathname	18
index	206
input file descriptor	
duplication	243
input mode	70
input processing	69
canonical mode	69
non-canonical mode	69
inter-user communication	270, 276
interactive facilities	269, 275
interface characteristics	67
interfaces	177
internationalization variable	56
interprocess communication	103
invalid	
use in RE	60
ioctl()	66, 91
IPC	103
ISO/IEC 646: 1991 standard	24
ISO C standard	10, 66, 83-85, 87-88, 91, 93, 100, 203
ISTRIP	70
job control	18-21, 23, 67-68, 93, 100, 271
job control, implementing applications	21
job control, implementing shells	19
job control, implementing systems	21
jobs	229, 275
join	275

kernel.....	22	lseek()	111-112, 118, 162, 203
kernel entity	166	lstat	271-272
kill	229, 275	lstat()	26, 228
kill()	93-95, 97-98, 100, 203	lutime()	26
last close	118	macro	8
lchown()	26	mail, rationale for omission	261
LC_COLLATE	48	mailx.....	275-276
LC_CTYPE.....	47	make.....	277-278
LC_MESSAGES	52	makecontext	206
LC_MONETARY	50	malloc().....	101, 121, 123, 151, 163-164, 173-174
LC_NUMERIC.....	51, 73	man	275
LC_TIME	51	map	22
ld, rationale for omission.....	261	mapped	22
legacy.....	7	margin code notation.....	9
legacy, rationale	7	margin codes	8, 83, 221
lex	277-278	mathematical functions	
library routine	22	error conditions	41
limits.....	221	NaN arguments	41
line, rationale for omission.....	261	MAX_CANON.....	69, 223, 285
LINES.....	57	MAX_INPUT.....	223, 286
LINE_MAX.....	23, 31, 222	may.....	6
link	272	may, rationale.....	6
link().....	16, 26	MCL_FUTURE.....	115
LINK_MAX	222, 285	MCL_INHERIT	116
lint, rationale for omission	261	memory locking.....	114
lio_listio().....	97, 112	memory management.....	113, 271
lists.....	247	memory management unit.....	114
ln.....	228, 276	memory object.....	22
local mode	71	memory synchronization	37
locale.....	45, 276	memory-resident	22
definition example	52	mesg	276-277
definition grammar.....	52	message passing	105, 271, 273
grammar	52	message queue.....	105
lexical conventions.....	52	mkdir	272, 276
locale configuration.....	269, 276	mkfifo.....	276
locale definition	46	mkfifo().....	18
localedef.....	276-277	mknod().....	18
localtime()	40, 163	mknod, rationale for omission	261
logger.....	277	mktemp	206
logical device.....	22	mktime().....	40
login, rationale for omission	261	mlockall()	115
LOGIN_NAME_MAX	285	mmap()	117-119, 121
LOGNAME.....	5, 57	MMU	114
logname	277	modem disconnect	70
longjmp()	90, 101, 171-172, 274	monotonic clock.....	137
LONG_MAX	72	more	275-276
LONG_MIN	72	mount point.....	22
lorder, rationale for omission.....	261	mount().....	22
lp.....	277	mounted file system.....	22
lpstat, rationale for omission	261	mprotect().....	118
ls.....	228, 276	mq_open().....	106
lseek.....	272	MQ_OPEN_MAX	286
		MQ_PRIO_MAX.....	286

Index

mq_receive()	107
mq_send()	107
mq_timedreceive()	138
mq_timedsend()	138
msg*()	103
msgctl()	103
msgget()	103
msgrcv()	103
msgsnd()	103
msync()	118
multi-byte character	17, 69-70
multiple tasks	269, 275
munmap()	117-119, 121, 124
mutex	155
mutex attributes	
extended	158
mutex initialization	170
mv	228, 276
name space	85
name space pollution	84-85
NAME_MAX	223, 286
NaN arguments	41
nanosleep()	135, 137-138, 273
native implementation	23
network byte order	36
newgrp	277
news, rationale for omission	262
NGROUPS_MAX	5, 25, 223, 280, 286
nice	275
nice value	23
nice()	128
nl_langinfo()	274
nm	277
noclobber option	241
nohup	275
non-canonical mode input processing	69
non-printable	31
normative references	6, 79, 219
NQS	255
obsolescent	6
obsolescent, rationale	6
od	275, 277
open	272
open file description	23
open file descriptors	243
open()	18, 68, 118-121, 228
opendir	272
openlog()	277
OPEN_MAX	5, 223, 286-287
option definitions	255
optional behavior	287
options	178
OR lists	248
orphaned process group	23, 100
output device	66
output file descriptor	
duplication	243
output mode	70
output processing	70
overflow conditions	202
overflow in dumping trace streams	202
overflow in trace streams	202
pack, rationale for omission	262
page	24, 117, 121
PAGESIZE	117, 162, 286
parallel I/O	162
parameter expansion	237
parameters	70, 233
parameters, positional	233
parameters, special	233
parent directory	24
passwd file	24
passwd, rationale for omission	262
paste	275
patch	276-277
PATH	57
pathchk	276
pathconf	272
pathconf()	18, 221, 271
pathname expansion	241
pathname resolution	38
pathname, incomplete	18
PATH_MAX	90, 223, 286
pattern matching	
multiple character	253
notation	252
single character	252
patterns	
filename expansion	253
pause	271
pause()	99, 102
pax	276-277
pcat, rationale for omission	262
pclose()	277
Pending error	177
per-thread errno	91
performance enhancements	268
pg, rationale for omission	262
PID_MAX	203
pipe	19, 24, 272
pipe()	99
pipelines	246
PIPE_BUF	223, 286
pointer types	204
popen()	274, 277-278

portability	8, 83, 221
portability codes	8, 83, 221
portable character set	42
portable filename character set	24
positional parameters	233
POSIX locale	46
POSIX.1 symbols	84
POSIX.13	121
POSIX2_BC_BASE_MAX	280
POSIX2_BC_DIM_MAX	280
POSIX2_BC_SCALE_MAX	280
POSIX2_BC_STRING_MAX	280
POSIX2_CHAR_TERM	279
POSIX2_COLL_WEIGHTS_MAX	280
POSIX2_C_BIND	223, 278
POSIX2_C_DEV	223, 278
POSIX2_EXPR_NEST_MAX	280
POSIX2_FORT_DEV	223, 279
POSIX2_FORT_RUN	223, 279
POSIX2_LINE_MAX	280
POSIX2_LOCALEDEF	224, 276, 279
POSIX2_PBS	279
POSIX2_PBS_ACCOUNTING	279
POSIX2_PBS_CHECKPOINT	279
POSIX2_PBS_LOCATE	279
POSIX2_PBS_MESSAGE	279
POSIX2_PBS_TRACK	279
POSIX2_RE_DUP_MAX	280
POSIX2_SW_DEV	223, 278
POSIX2_SYMLINKS	223
POSIX2_UPE	223, 278-279
POSIX2_VERSION	280
POSIX_ALLOC_SIZE_MIN	104
POSIX_ASYNCHRONOUS_IO	291
POSIX_BARRIERS	291
POSIX_CLOCK_SELECTION	292
POSIX_C_LANG_JUMP	291
POSIX_C_LANG_MATH	291
POSIX_C_LANG_SUPPORT	292
POSIX_C_LANG_SUPPORT_R	292
POSIX_C_LANG_WIDE_CHAR	292
POSIX_C_LANG_WIDE_CHAR_EXT	292
POSIX_C_LIB_EXT	292
POSIX_DEVICE_IO	292
POSIX_DEVICE_IO_EXT	292
POSIX_DEVICE_SPECIFIC	292
POSIX_DEVICE_SPECIFIC_R	293
POSIX_DYNAMIC_LINKING	293
posix_fadvise()	104
POSIX_FADV_DONTNEED	104
POSIX_FADV_NOREUSE	104
POSIX_FADV_RANDOM	104
POSIX_FADV_SEQUENTIAL	104
POSIX_FADV_WILLNEED	104
POSIX_FD_MGMT	293
POSIX_FIFO	293
POSIX_FIFO_FD	293
POSIX_FILE_ATTRIBUTES	293
POSIX_FILE_ATTRIBUTES_FD	293
POSIX_FILE_LOCKING	293
POSIX_FILE_SYSTEM	293
POSIX_FILE_SYSTEM_EXT	293
POSIX_FILE_SYSTEM_FD	293
POSIX_FILE_SYSTEM_GLOB	293
POSIX_FILE_SYSTEM_R	293
POSIX_I18N	293
POSIX_JOB_CONTROL	293
posix_madvise()	104
POSIX_MADV_DONTNEED	104
POSIX_MADV_RANDOM	104
POSIX_MADV_SEQUENTIAL	104
POSIX_MADV_WILLNEED	104
POSIX_MAPPED_FILES	293
POSIX_MEMORY_PROTECTION	293
posix_mem_offset()	121-122
POSIX_MULTI_CONCURRENT_LOCALES	293
POSIX_MULTI_PROCESS	294
POSIX_MULTI_PROCESS_FD	294
POSIX_NETWORKING	294
POSIX_PIPE	294
POSIX_REALTIME_SIGNALS	294
POSIX_REC_INCR_XFER_SIZE	105
POSIX_REC_MAX_XFER_SIZE	105
POSIX_REC_MIN_XFER_SIZE	105
POSIX_REC_XFER_ALIGN	104
POSIX_REGEX	294
POSIX_ROBUST_MUTEXES	294
POSIX_RW_LOCKS	294
POSIX_SEMAPHORES	294
POSIX_SHELL_FUNC	294
POSIX_SIGNALS	294
POSIX_SIGNALS_EXT	294
POSIX_SIGNAL_JUMP	294
POSIX_SINGLE_PROCESS	295
posix_spawn()	207, 271
posix_spawnp()	207, 271
POSIX_SPIN_LOCKS	295
POSIX_SYMBOLIC_LINKS	295
POSIX_SYMBOLIC_LINKS_FD	295
POSIX_SYSTEM_DATABASE	295
POSIX_SYSTEM_DATABASE_R	295
POSIX_THREADS_BASE	295
POSIX_THREADS_EXT	295
POSIX_TIMERS	295
posix_trace_eventid_open()	197

Index

POSIX_TRACE_LOOP	202
posix_typed_mem_get_info()	121
posix_typed_mem_open()	121
POSIX_USER_GROUPS	295
POSIX_USER_GROUPS_R.....	295
POSIX_VERSION	286
POSIX_WIDE_CHAR_DEVICE_IO.....	295
post-mortem filtering of trace event types	199
pr.....	275, 277
pread()	163
printf	274-275
printing	270
privilege	34
process group	67
process group ID	19, 67
process group lifetime	67
process group, orphaned.....	23, 100
process groups, concepts in job control.....	19
process ID reuse.....	39
process ID, rationale.....	203
process lifetime	25
process management.....	268, 271
process scheduling	127, 271
process termination.....	25
prof, rationale for omission.....	262
profiling	278
programming manipulation	193
prompting.....	235
protocols	177
ps.....	275, 277
pthread	200
pthread_attr_getguardsize().....	162
pthread_attr_getstackaddr	206
pthread_attr_setguardsize()	162
pthread_attr_setstackaddr	206
PTHREAD_BARRIER_SERIAL_THREAD.....	153
pthread_barrier_wait().....	154, 173
pthread_cond_init().....	150
pthread_cond_timedwait().....	92, 137, 159, 281
pthread_cond_wait()	92, 109, 159
pthread_create().....	150-151
PTHREAD_CREATE_DETACHED	171
PTHREAD_DESTRUCTOR_ITERATIONS	286
pthread_detach().....	171
pthread_getconcurrency()	162
pthread_getcpuclockid()	140-141
pthread_join()	92, 171
PTHREAD_KEYS_MAX.....	286
pthread_key_create().....	152
pthread_mutexattr_gettype().....	160
pthread_mutexattr_settype()	160
PTHREAD_MUTEX_DEFAULT.....	159
PTHREAD_MUTEX_ERRORCHECK.....	159
pthread_mutex_init()	150
pthread_mutex_lock().....	92, 159, 173
PTHREAD_MUTEX_NORMAL	159
PTHREAD_MUTEX_RECURSIVE	159
pthread_mutex_timedlock().....	138
pthread_mutex_trylock().....	159
pthread_mutex_unlock()	159
PTHREAD_PROCESS_PRIVATE	161
PTHREAD_PROCESS_SHARED	161
pthread_rwlockattr_destroy().....	161
pthread_rwlockattr_getpshared()	161
pthread_rwlockattr_init()	160
pthread_rwlockattr_setpshared().....	161
pthread_rwlock_init()	161
PTHREAD_RWLOCK_INITIALIZER.....	161
pthread_rwlock_rdlock().....	161
pthread_rwlock_t	160
pthread_rwlock_tryrdlock().....	161
pthread_rwlock_trywrlock()	161
pthread_rwlock_unlock()	161, 175
pthread_rwlock_wrlock()	161
pthread_self()	152
pthread_setconcurrency().....	162
pthread_setprio()	170
pthread_setschedparam()	170
pthread_setspecific()	152
pthread_spin_lock().....	154, 173
pthread_spin_trylock().....	154
PTHREAD_STACK_MIN.....	286
PTHREAD_THREADS_MAX.....	286
putc	272
putc().....	164
putchar	272
pwd	276
pwrite().....	163
queuing of waiting threads.....	175
quote removal	241
quoting.....	230
rand()	173
RCS, rationale for omission.....	262
RE	
grammar	65
RE bracket expression.....	61
read	229, 272, 274
read lock.....	160
read-write attributes	160
read-write locks	160
readdir.....	272
reading an active trace stream	202
reading data	69
readlink	272

readlink()	26
realpath	272
realtime	104
realtime signal delivery	96
realtime signal generation	96
realtime signals	109
red, rationale for omission	262
redirect input	243
redirect output	243
redirection	241
references	6, 79, 219
regcomp()	277
regerror()	277
regexec()	277
regfree()	277
regular expression	58
definitions	59
general requirements	60
grammar	64
regular file	25
rejected utilities	260
remove()	26
rename	272
rename()	26
renice	275
replenishment period	130
reserved words	233
rewinddir	272
RE_DUP_MAX	223
rindex	206
rm	228, 276
rmdir	272, 276
rmdir()	26, 91
robust mutexes	156
root directory	25, 39
root file system	25
root of a file system	25
routing	177
rsh, rationale for omission	262
RTSIG_MAX	287
samefile()	203
SA_NOCLDSTOP	20
SA_RESETHAND	205
SA_RESTART	205
SA_SIGINFO	97-98
scalb	207
scheduling allocation domain	168
scheduling contention scope	168-169
scheduling documentation	168
scheduling policy	40
SCHED_FIFO	128-130, 168, 175, 273
SCHED_OTHER	129
SCHED_RR	128-129, 168, 175, 273
SCHED_SPORADIC	273
scope	3, 79, 219
sdb, rationale for omission	262
sdiff, rationale for omission	262
seconds since the Epoch	40
security considerations	13, 17, 21, 32, 34, 68
security, monolithic privileges	13
sed	275-276
sem*()	103
semaphore	41, 107, 273
semctl()	103
semget()	103
semop()	103
sem_init()	107
SEM_NSEMS_MAX	287
sem_open()	107
sem_timedwait()	138
sem_trywait()	92, 109
SEM_VALUE_MAX	287
sem_wait()	92, 109
sequential lists	248
session	20, 24, 68
set	235
setgid()	25
setgrent()	32
setjmp()	274
setlocale()	274
setlogmask()	277
setpgid()	19-20, 67
setpriority()	128
setpwent()	32
setrlimit()	228
setsid()	67
setuid()	25
sh	276, 283
shall	6
shall, rationale	6
shar, rationale for omission	262
shared memory	119
shell	18-21
SHELL	58
shell	67-68, 93, 99-100
shell commands	244
shell errors	243
shell execution environment	233, 252
shell grammar	251
lexical conventions	252
rules	252
shell variables	234
shell, job control	19, 93, 100
shl, rationale for omission	262
shm*()	103

Index

shmctl()	103
shmdt()	103
shm_open()	118-121
shm_unlink()	119-121
should	6
should, rationale	6
SIGABRT	30, 93
sigaction()	95, 97
SIGBUS	30, 93
SIGCHLD	20, 95, 99
SIGCLD	99
SIGCONT	20, 95, 99-100
SIGEMT	93
SIGEV_NONE	96
SIGEV_SIGNAL	96
SIGEV_THREAD	96
SIGFPE	30, 93, 95, 98
SIGHUP	100
SIGILL	30, 93
SIGINT	21, 166
SIGIOT	93
SIGKILL	93, 95, 100
siglongjmp()	90, 101, 274
signal	25
signal acceptance	94
signal actions	99
signal concepts	93
signal delivery	94
signal generation	94
signal names	93
signal()	93, 95
signals	178, 252
SIGPIPE	30, 91
sigprocmask()	95
SIGRTMAX	97-98
SIGRTMIN	97-98
SIGSEGV	30, 93, 98
sigsetjmp()	274
sigset_t	93
SIGSTOP	99
sigsuspend()	99, 102
SIGSYS	93
SIGTERM	93
sigtimedwait()	92, 111, 138
SIGTRAP	93
SIGTSTP	21, 99-100
SIGTTIN	20, 68, 99-100
SIGTTOU	20, 68, 99-100
SIGUSR1	93
SIGUSR2	93
sigwait()	92, 172
sigwaitinfo()	92, 111
sigwait_multiple()	95
SIG_DFL	95-96, 99
SIG_IGN	20, 95-96, 99, 101
simple commands	244
single-quotes	230
size, rationale for omission	262
SI_USER	97-98
sleep	271, 274
sleep()	101-102, 273
socket I/O mode	177
socket out-of-band data state	178
socket owner	177
socket queue limit	177
socket receive queue	178
socket types	177
sockets	177
Internet Protocols	178
IPv4	178
IPv6	178
local UNIX connection	178
software development	270, 277
sort	275-276
spawn example	207
special built-in	249
special built-in utilities	254
special characters	70
special control character	71
special parameters	233
specific implementation	18
spell, rationale for omission	263
spin locks	154-155
split	275
sporadic server scheduling policy	130
SSIZE_MAX	204, 287
SS_REPL_MAX	133
standard I/O streams	102
stat	228, 271-272
stat()	16, 118, 228
state-dependent character encoding	43
statvfs()	228
STREAMS	102
STREAM_MAX	287
strings	277
strip	277
structures, additions to	85
stty	57, 277
su, rationale for omission	263
subprofiling	11
subprofiling option groups	291
subshells	20
successfully completed	31
sum	228
sum, rationale for omission	263

superuser	13, 25, 34, 235, 260
supplementary group ID	25
supplementary groups	34
Supported Threads functions	157
swapcontext	206
symbolic constant	8
symbolic link	26
symbolic name	8
symbols	84
SYMLOOP_MAX	90
synchronized I/O	112, 272
data integrity completion	31, 112
file integrity completion	31, 112
synchronously-generated signal	30
sysconf()	18, 115, 117, 119, 166, 221, 271-272
syslog()	277
system call	31
system console	31
System database	31
system documentation	7
system environment	270, 277
System III	24, 203
system interfaces	205
system process	31
system reboot	31
System V	16, 21, 93, 95, 99
system()	274, 277-278
tabs	275
tail	275
talk	276
tar, rationale for omission	263
tcgetattr()	20
tcgetpgrp()	20, 67
tcsetattr()	20, 66
tcsetpgrp()	19-20
tee	274
TERM	57
terminal access control	68
terminal device file	67
closing	70
terminal type	66
terminology	6, 82, 220
termios structure	70
test	274, 276
TeX	276
text file	31
thread	32
thread cancelability states	173
thread cancelability type	173
thread cancellation	171, 173
thread cancellation points	173
thread concurrency level	161
thread creation attributes	149
thread ID	32, 166
thread interactions	176
thread mutex	166
thread read-write lock	175
thread scheduling	166
thread stack guard size	162
thread-safe	32
thread-safe function	32
thread-safety	41, 163
thread-safety, rationale	41
thread-specific data	151
threads	149
implementation models	151
threads extensions	158
tilde expansion	236
time	274, 277
time()	90
timeouts	142
timers	133
TIMER_ABSTIME	134-136
TIMER_MAX	287
timer_settime()	134-136
times()	90, 141, 271
timestamp clock	201
time_t	40
token recognition	232
TOSTOP	20
touch	228, 276
tput	277
tr	275-276
trace analyzer	191
trace event type-filtering	199
trace event types	199
trace examples	188
trace model	183
trace operation control	188
trace storage	187
trace stream attribute	193
trace stream states	186
tracing	41, 178
tracing all processes	187
tracing, detailed objectives	179
triggering	201
troff	276
trojan horse	13
true	229, 274
tty	277
ttyname()	163
TTY_NAME_MAX	287
typed memory	121
TZ	58
TZNAME_MAX	287

Index

tzset()	274
ualarm	207, 271
UID_MAX	204
uid_t	32
ULONG_MAX	221
umask	229, 277
umask()	271
umount()	22
unalias	229, 275
uname	277
uname()	271
unbounded priority inversion	170
undefined	7
undefined, rationale	7
unexpand	275
uniq	275-276
unlink	272
unlink()	26, 91, 118-119, 121
unpack, rationale for omission	263
unsafe functions	100
unspecified	7
unspecified, rationale	7
until loop	249
user database	32
user database access	33
user requirements	267
usleep	207, 271
utility	41
utility argument syntax	71
utility conventions	71
utility description defaults	224
utility limits	221
utility syntax guidelines	72
utime	272
utime()	26
uucp	276
uudecode	275-276
uuencode	275-276
variable assignment	41
variables	233
VEOF	71
VEOL	71
Version 7	39, 230
vfork	207
vhangup()	21
vi	275-276
virtual processor	33
VMIN	71
VTIME	71
wait	229, 274
wait()	90, 93, 99-100, 102, 271
waitid()	271
waitpid()	20, 24, 100, 203, 271
wall, rationale for omission	263
wc	275
wcswcs	207
WERASE	69
while loop	249
who	276-277
wide-character codes	43
word expansions	236
wordexp()	277
wordfree()	277
write	272, 275-276
write lock	160
writes	20, 68, 90, 99-100, 102, 111-113, 117-118, 204
writing data	70
WUNTRACED	20
xargs	274
XSI	33
XSI IPC	103
XSI_C_LANG_SUPPORT	295
XSI_DBM	296
XSI_DEVICE_IO	296
XSI_DEVICE_SPECIFIC	296
XSI_FILE_SYSTEM	296
XSI_IPC	296
XSI_JUMP	296
XSI_MATH	296
XSI_MULTI_PROCESS	296
XSI_SIGNALS	296
XSI_SINGLE_PROCESS	296
XSI_SYSTEM_DATABASE	296
XSI_SYSTEM_LOGGING	296
XSI_THREADS_EXT	296
XSI_TIMERS	296
XSI_USER_GROUPS	296
XSI_WIDE_CHAR	296
yacc	277-278

