

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

Prepared by the Austin Group
(<http://www.opengroup.org/austin/>)

Copyright © 2001 The Institute of Electrical and Electronics Engineers, Inc.
3 Park Avenue, New York, NY 10016-5997, USA.

Copyright © 2001 The Open Group
Apex Plaza, Forbury Road, Reading, Berkshire RG11AX, 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.

1 / *Rationale (Informative)* |

2 **Part A:** |
3 **Base Definitions** |

4 *The Open Group* |
5 *The Institute of Electrical and Electronics Engineers, Inc.* |

Rationale for Base Definitions

7

8 A.1 Introduction

9 A.1.1 Scope

10 IEEE Std 1003.1-200x is one of a family of standards known as POSIX. The family of standards
11 extends to many topics; IEEE Std 1003.1-200x is known as POSIX.1 and consists of both
12 operating system interfaces and shell and utilities. IEEE Std 1003.1-200x is technically identical
13 to The Open Group Base Specifications, Issue 6, which comprise the core volumes of the Single
14 UNIX Specification, Version 3.

15 Scope of IEEE Std 1003.1-200x

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

19 Since the revision includes merging the Base volumes of the Single UNIX Specification, many
20 features that were previously not *adopted* into earlier revisions of POSIX.1 and POSIX.2 are now
21 included in IEEE Std 1003.1-200x. In most cases, these additions are part of the XSI extension; in
22 other cases the standard developers decided that now was the time to migrate these to the base
23 standard.

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

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

31 The following IEEE Standards have been added to the base documents in this revision:

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

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

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

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

- 49 • Backward-compatibility

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

- 54 • Architecture and n-bit neutral

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

- 58 • Extensibility

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

62 **POSIX.1 and the ISO C standard**

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

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

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

71 The following block of text was added to each reference page affected:

72 *The functionality described on this reference page is aligned with the ISO C standard. Any conflict*
 73 *between the requirements described here and the ISO C standard is unintentional. This volume of*
 74 *IEEE Std 1003.1-200x defers to the ISO C standard.*

75 and each page was parsed for additions beyond the ISO C standard (that is, including both
 76 POSIX and UNIX extensions), and these extensions are marked as CX extensions (for C
 77 Extensions). |

FIPS Requirements

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

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

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

A.1.2 Conformance

See Section A.2 (on page 3299).

A.1.3 Normative References

There is no additional rationale provided for this section.

A.1.4 Terminology

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

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

conformance document

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

116 implementation-defined

117 This definition is analogous to that of the ISO C standard and, together with *undefined* and
118 *unspecified*, provides a range of specification of freedom allowed to the interface
119 implementor.

120 may

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

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

126 shall

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

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

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

132 are meant to be exactly equivalent.

133 should

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

137 obsolescent

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

145 legacy

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

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

159 system documentation

160 The system documentation should normally describe the whole of the implementation,
161 including any extensions provided by the implementation. Such documents normally
162 contain information at least as detailed as the specifications in IEEE Std 1003.1-200x. Few

163 requirements are made on the system documentation, but the term is needed to avoid a
 164 dangling pointer where the conformance document is permitted to point to the system
 165 documentation.

166 **undefined**

167 See *implementation-defined*.

168 **unspecified**

169 See *implementation-defined*.

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

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

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

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

195 In many places IEEE Std 1003.1-200x is silent about the behavior of some possible construct.
 196 For example, a variable may be defined for a specified range of values and behaviors are
 197 described for those values; nothing is said about what happens if the variable has any other
 198 value. That kind of silence can imply an error in the standard, but it may also imply that the
 199 standard was intentionally silent and that any behavior is permitted. There is a natural
 200 tendency to infer that if the standard is silent, a behavior is prohibited. That is not the intent.
 201 Silence is intended to be equivalent to the term *unspecified*.

202 The term *application* is not defined in IEEE Std 1003.1-200x; it is assumed to be a part of
 203 general computer science terminology.

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

206 **macro** |

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

210 **symbolic name** |
 211 This can also refer to a C preprocessor symbol (without arguments), but is also used to refer |
 212 to the names for characters in character sets. In addition, it is sometimes used to refer to |
 213 host names and even filenames. |

214 **symbolic constant** |
 215 This also refers to a C preprocessor symbol (also without arguments). |

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

218 A.1.5 Portability

219 To aid the identification of options within IEEE Std 1003.1-200x, a notation consisting of margin |
 220 codes and shading is used. This is based on the notation used in previous revisions of The Open |
 221 Group's Base specifications.

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

227 A.1.5.1 Codes

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

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

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

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

241 OB This margin code is used to denote obsolescent behavior and thus flag a possible future |
 242 application portability warning.

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

247 XSI This code is used to denote interfaces and facilities within interfaces only required on |
 248 systems supporting the XSI extension. This is introduced to support the Single UNIX |
 249 Specification.

250 XSR This code is used to denote interfaces and facilities within interfaces only required on |
 251 systems supporting STREAMS. This is introduced to support the Single UNIX |
 252 Specification, although it is defined in a way so that it can standalone from the XSI |
 253 notation.

254 A.1.5.2 *Margin Code Notation*

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

260 **A.2 Conformance**

261 The terms *profile* and *profiling* are used throughout this section.

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

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

- 266 • The addition of profiling options, allowing larger profiles of options such as the XSI |
 267 extension used by the Single UNIX Specification. In effect, it has profiled itself (that is, |
 268 created a self-profile). |
- 269 • The addition of a definition of subprofiling considerations, to allow smaller profiles of |
 270 options. |
- 271 • The addition of a hierarchy of super-options for XSI; these were formerly known as *Feature* |
 272 *Groups* in The Open Group System Interfaces and Headers, Issue 5 specification. |
- 273 • Options from the ISO POSIX-2: 1993 standard are also now included as IEEE Std 1003.1-200x |
 274 merges the functionality from it.

275 **A.2.1 Implementation Conformance**

276 These definitions allow application developers to know what to depend on in an
 277 implementation.

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

282 A.2.1.1 *Requirements*

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

286 A.2.1.2 *Documentation*

287 The conformance documentation is required to use the same numbering scheme as POSIX.1 for
 288 purposes of cross-referencing. All options that an implementation chooses shall be reflected in
 289 <limits.h> and <unistd.h>.

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

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

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

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

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

319 A.2.1.3 *POSIX Conformance*

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

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

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

- 329 • `_POSIX_JOB_CONTROL`
- 330 • `_POSIX_SAVED_IDS`
- 331 • `_POSIX_VDISABLE`

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

336 A.2.1.4 *XSI Conformance*

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

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

341 A.2.1.5 *Option Groups*

342 The concept of *Option Groups* is introduced to IEEE Std 1003.1-200x to allow collections of
343 related functions or options to be grouped together. This has been used as follows: the *XSI*
344 *Option Groups* have been created to allow super-options, collections of underlying options and
345 related functions, to be collectively supported by XSI-conforming systems. These reflect the
346 *Feature Groups* from The Open Group System Interfaces and Headers, Issue 5 specification.

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

351 A.2.1.6 *Options*

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

355 **A.2.2 Application Conformance**

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

359 POSIX.1 occasionally uses the expressions *portable application* or *conforming application*. As they
360 are used, these are synonyms for any of these terms. The differences between the classes of
361 application conformance relate to the requirements for other standards, the options supported
362 (such as the XSI extension) or, in the case of the Conforming POSIX.1 Application Using
363 Extensions, to implementation extensions. When one of the less explicit expressions is used, it
364 should be apparent from the context of the discussion which of the more explicit names is
365 appropriate

366 A.2.2.1 *Strictly Conforming POSIX Application*

367 This definition is analogous to that of a ISO C standard *conforming program*.

368 The major difference between a *Strictly Conforming POSIX Application* and a ISO C standard
369 *strictly conforming program* is that the latter is not allowed to use features of POSIX that are not in
370 the ISO C standard.

371 A.2.2.2 *Conforming POSIX Application*

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

373 A.2.2.3 *Conforming POSIX Application Using Extensions*

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

378 A.2.2.4 *Strictly Conforming XSI Application*

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

381 A.2.2.5 *Conforming XSI Application Using Extensions*

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

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

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

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

391 There is no additional rationale provided for this section.

392 **A.3 Definitions**

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

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

401 Many of these definitions are necessarily circular, and some of the terms (such as *process*) are |
 402 variants of basic computing science terms that are inherently hard to define. Where some |
 403 definitions are more conceptual and contain requirements, these appear in the General Concepts |
 404 chapter. Those listed in this section appear in an alphabetical glossary format of terms.

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

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

413 **Appropriate Privileges**

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

428 **Byte**

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

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

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

440 The standard developers also felt that this was not an undue restriction for the current state of
 441 the art for this version of IEEE Std 1003.1-200x, but recognize that if industry trends continue, a
 442 wider character type may be required in the future.

443 **Character**

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

448 **Clock Tick**

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

452 **Command**

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

465 The terms *tool* and *application* can be thought of as being synonymous with *utility* from the
466 perspective of the operating system kernel. Tools, applications, and utilities historically have
467 run, typically, in processes above the kernel level. Tools and utilities historically have been a part
468 of the operating system non-kernel code and have performed system-related functions, such as
469 listing directory contents, checking file systems, repairing file systems, or extracting system
470 status information. Applications have not generally been a part of the operating system, and
471 they perform non-system-related functions, such as word processing, architectural design,
472 mechanical design, workstation publishing, or financial analysis. Utilities have most frequently
473 been provided by the operating system distributor, applications by third-party software
474 distributors, or by the users themselves. Nevertheless, IEEE Std 1003.1-200x does not
475 differentiate between tools, utilities, and applications when it comes to receiving services from
476 the system, a shell, or the standard utilities. (For example, the *xargs* utility invokes another
477 utility; it would be of fairly limited usefulness if the users could not run their own applications
478 in place of the standard utilities.) Utilities are not applications in the sense that they are not
479 themselves subject to the restrictions of IEEE Std 1003.1-200x or any other standard—there is no
480 requirement for *grep*, *stty*, or any of the utilities defined here to be any of the classes of
481 conforming applications.

482 **Column Positions**

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

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

498 The term *column position* was used instead of the more natural *column* because the latter is
499 frequently used in the different contexts of columns of figures, columns of table values, and so
500 on. Wherever confusion might result, these latter types of columns are referred to as *text*

501 *columns.*

502 **Controlling Terminal**

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

506 **Device Number***

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

508 **Direct I/O**

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

511 **Directory**

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

516 **Directory Entry**

517 Throughout IEEE Std 1003.1-200x, the term *link* is used (about the *link()* function, for example)
518 in describing the objects that point to files from directories.

519 **Display**

520 The Shell and Utilities volume of IEEE Std 1003.1-200x assigns precise requirements for the
521 terms *display* and *write*. Some historical systems have chosen to implement certain utilities
522 without using the traditional file descriptor model. For example, the *vi* editor might employ
523 direct screen memory updates on a personal computer, rather than a *write()* system call. An
524 instance of user prompting might appear in a dialog box, rather than with standard error. When
525 the Shell and Utilities volume of IEEE Std 1003.1-200x uses the term *display*, the method of
526 outputting to the terminal is unspecified; many historical implementations use *termcap* or
527 *terminfo*, but this is not a requirement. The term *write* is used when the Shell and Utilities volume
528 of IEEE Std 1003.1-200x mandates that a file descriptor be used and that the output can be
529 redirected. However, it is assumed that when the writing is directly to the terminal (it has not
530 been redirected elsewhere), there is no practical way for a user or test suite to determine whether
531 a file descriptor is being used. Therefore, the use of a file descriptor is mandated only for the
532 redirection case and the implementation is free to use any method when the output is not
533 redirected. The verb *write* is used almost exclusively, with the very few exceptions of those
534 utilities where output redirection need not be supported: *tabs*, *talk*, *tput*, and *vi*.

535 Dot

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

538 Dot-Dot

539 Historical implementations permit the use of these filenames without their special meanings.
540 Such use precludes any meaningful use of these filenames by a Conforming POSIX.1
541 Application. Therefore, such use is considered an extension, the use of which makes an
542 implementation non-conforming; see also Section A.4.11 (on page 3327).

543 Epoch

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

548 FIFO Special File

549 See *pipe* in **Pipe** (on page 3314).

550 File

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

552 File Classes

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

556 Filename

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

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

566 File System

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

572 **Graphic Character**

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

575 **Group Database**

576 See **User Database** (on page 3322).

577 **Group File***

578 Implementation-defined; see **User Database** (on page 3322).

579 **Historical Implementations***

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

582 **Hosted Implementation***

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

590 **Implementation***

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

594 The term *specific implementation* is sometimes used as a synonym for *implementation*. This should
595 not be interpreted too narrowly; both terms can represent a relatively broad group of systems.
596 For example, a hardware vendor could market a very wide selection of systems that all used the
597 same instruction set, with some systems desktop models and others large multi-user
598 minicomputers. This wide range would probably share a common POSIX.1 operating system,
599 allowing an application compiled for one to be used on any of the others; this is a
600 [*specific*]implementation.

601 However, that wide range of machines probably has some differences between the models.
602 Some may have different clock rates, different file systems, different resource limits, different
603 network connections, and so on, depending on their sizes or intended usages. Even on two
604 identical machines, the system administrators may configure them differently. Each of these
605 different systems is known by the term a *specific instance of a specific implementation*. This term is
606 only used in the portions of POSIX.1 dealing with runtime queries: *sysconf()* and *pathconf()*.

607 Incomplete Pathname*

608 Absolute pathname has been adequately defined.

609 Job Control

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

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

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

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

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

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

636 Implementing Job Control Shells

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

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

643 The concept of *process groups* is used to provide this grouping. The shell places each job in a
644 separate process group via the *setpgid()* function. To do this, the *setpgid()* function is invoked by
645 the shell for each process in the job. It is actually useful to invoke *setpgid()* twice for each
646 process: once in the child process, after calling *fork()* to create the process, but before calling one
647 of the *exec* family of functions to begin execution of the program, and once in the parent shell
648 process, after calling *fork()* to create the child. The redundant invocation avoids a race condition
649 by ensuring that the child process is placed into the new process group before either the parent
650 or the child relies on this being the case. The *process group ID* for the job is selected by the shell to
651 be equal to the *process ID* of one of the processes in the job. Some shells choose to make one
652 process in the job be the parent of the other processes in the job (if any). Other shells (for

653 example, the C Shell) choose to make themselves the parent of all processes in the pipeline (job).
654 In order to support this latter case, the *setpgid()* function accepts a process group ID parameter
655 since the correct process group ID cannot be inherited from the shell. The shell itself is
656 considered to be a job and is the sole process in its own process group.

657 The shell also controls which job is currently in the foreground. A foreground and background
658 job differ in two ways: the shell waits for a foreground command to complete (or stop) before
659 continuing to read new commands, and the terminal I/O driver inhibits terminal access by
660 background jobs (causing the processes to stop). Thus, the shell must work cooperatively with
661 the terminal I/O driver and have a common understanding of which job is currently in the
662 foreground. It is the user who decides which command should be currently in the foreground,
663 and the user informs the shell via shell commands. The shell, in turn, informs the terminal I/O
664 driver via the *tcsetpgrp()* function. This indicates to the terminal I/O driver the process group ID
665 of the foreground process group (job). When the current foreground job either stops or
666 terminates, the shell places itself in the foreground via *tcsetpgrp()* before prompting for
667 additional commands. Note that when a job is created the new process group begins as a
668 background process group. It requires an explicit act of the shell via *tcsetpgrp()* to move a
669 process group (job) into the foreground.

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

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

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

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

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

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

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

698 A subshell can be spawned to perform an interactive function (prompting the terminal for
699 commands) or a non-interactive function (reading commands from a file). When operating non-
700 interactively, the job control shell will refrain from performing the job control-specific actions

701 described above. It will behave as a shell that does not support job control. For example, all *jobs*
702 will be left in the same process group as the shell, which itself remains in the process group
703 established for it by its parent. This allows the shell and its children to be treated as a single job
704 by a parent shell, and they can be affected as a unit by terminal keyboard signals.

705 An interactive subshell can be spawned from another job control-cognizant shell in either the
706 foreground or background. (For example, from the C Shell, the user can execute the command,
707 `csch &`.) Before the subshell activates job control by calling `setpgid()` to place itself in its own
708 process group and `tcsetpgrp()` to place its new process group in the foreground, it needs to
709 ensure that it has already been placed in the foreground by its parent. (Otherwise, there could
710 be multiple job control shells that simultaneously attempt to control mediation of the terminal.)
711 To determine this, the shell retrieves its own process group via `getpgrp()` and the process group
712 of the current foreground job via `tcgetpgrp()`. If these are not equal, the shell sends SIGTTIN to
713 its own process group, causing itself to stop. When continued later by its parent, the shell
714 repeats the process group check. When the process groups finally match, the shell is in the
715 foreground and it can proceed to take control. After this point, the shell ignores all the job
716 control stop signals so that it does not inadvertently stop itself.

717 *Implementing Job Control Applications*

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

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

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

743 *Implementing Job Control Systems*

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

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

750 POSIX.1 does not specify how controlling terminal access is affected by a user logging out (that
751 is, by a controlling process terminating). 4.2 BSD uses the *vhangup()* function to prevent any
752 access to the controlling terminal through file descriptors opened prior to logout. System V does
753 not prevent controlling terminal access through file descriptors opened prior to logout (except
754 for the case of the special file, */dev/tty*). Some implementations choose to make processes
755 immune from job control after logout (that is, such processes are always treated as if in the
756 foreground); other implementations continue to enforce foreground/background checks after
757 logout. Therefore, a Conforming POSIX.1 Application should not attempt to access the
758 controlling terminal after logout since such access is unreliable. If an implementation chooses to
759 deny access to a controlling terminal after its controlling process exits, POSIX.1 requires a certain
760 type of behavior (see **Controlling Terminal** (on page 3305)).

761 **Kernel***

762 See *system call*.

763 **Library Routine***

764 See *system call*.

765 **Logical Device***

766 Implementation-defined.

767 **Map**

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

770 **Memory-Resident**

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

784 Memory Object*

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

790 Mount Point*

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

793 Mounted File System*

794 See *file system*.

795 Name

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

807 Native Implementation*

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

812 Nice Value

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

820 Open File Description

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

824 The following alternate names were discussed:

- 825 • For *open file description*:
826 *open instance, file access description, open file information, and file access information.*
- 827 • For *file descriptor*:
828 *file handle, file number (c.f., `fileno()`). Some historical implementations use the term *file table*
829 *entry.**

830 Orphaned Process Group

831 Historical implementations have a concept of an orphaned process, which is a process whose
832 parent process has exited. When job control is in use, it is necessary to prevent processes from
833 being stopped in response to interactions with the terminal after they no longer are controlled by
834 a job control-cognizant program. Because signals generated by the terminal are sent to a process
835 group and not to individual processes, and because a signal may be provoked by a process that
836 is not orphaned, but sent to another process that is orphaned, it is necessary to define an
837 orphaned process group. The definition assumes that a process group will be manipulated as a
838 group and that the job control-cognizant process controlling the group is outside of the group
839 and is the parent of at least one process in the group (so that state changes may be reported via
840 *waitpid()*). Therefore, a group is considered to be controlled as long as at least one process in the
841 group has a parent that is outside of the process group, but within the session.

842 This definition of orphaned process groups ensures that a session leader’s process group is
843 always considered to be orphaned, and thus it is prevented from stopping in response to
844 terminal signals.

845 Page

846 The term *page* is defined to support the description of the behavior of memory mapping for
847 shared memory and memory mapped files, and the description of the behavior of process
848 memory locking. It is not intended to imply that shared memory/file mapping and memory
849 locking are applicable only to “paged” architectures. For the purposes of IEEE Std 1003.1-200x,
850 whatever the granularity on which an architecture supports mapping or locking is considered to
851 be a “page” . If an architecture cannot support the memory mapping or locking functions
852 specified by IEEE Std 1003.1-200x on any granularity, then these options will not be
853 implemented on the architecture.

854 Passwd File*

855 Implementation-defined; see **User Database** (on page 3322).

856 Parent Directory

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

863 Pipe

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

866 Portable Filename Character Set

867 The encoding of this character set is not specified—specifically, ASCII is not required. But the
868 implementation must provide a unique character code for each of the printable graphics
869 specified by POSIX.1; see also Section A.4.6 (on page 3324).

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

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

884 Regular File

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

888 Root Directory

889 This definition permits the operation of *chroot()*, even though that function is not in POSIX.1; see
890 also *file hierarchy*.

891 Root File System*

892 Implementation-defined.

893 Root of a File System*

894 Implementation-defined; see *mount point*.

895 Signal

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

898 Superuser*

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

901 Supplementary Group ID

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

- 909 1. Reword the description of *setgid()* to permit it to change the supplementary group IDs to
910 reflect the new effective group ID. A problem with this is that it adds more “may”s to the
911 wording and does not address the portability problems of this optional behavior.
- 912 2. Mandate the inclusion of the effective group ID in the supplementary set (giving
913 {NGROUPS_MAX} a minimum value of 1). This is the behavior of 4.4 BSD. In that system,
914 the effective group ID is the first element of the array of supplementary group IDs (there is
915 no separate copy stored, and changes to the effective group ID are made only in the
916 supplementary group set). By convention, the initial value of the effective group ID is
917 duplicated elsewhere in the array so that the initial value is not lost when executing a set-
918 group-ID program.
- 919 3. Change the definition of supplementary group ID to exclude the effective group ID and
920 specify that the effective group ID does not change the set of supplementary group IDs.
921 This is the behavior of 4.2 BSD, 4.3 BSD, and System V, Release 4.
- 922 4. Change the definition of supplementary group ID to exclude the effective group ID, and
923 require that *getgroups()* return the union of the effective group ID and the supplementary
924 group IDs.
- 925 5. Change the definition of {NGROUPS_MAX} to be one more than the number of
926 supplementary group IDs, so it continues to be the number of values returned by
927 *getgroups()* and existing applications continue to work. This alternative is effectively the
928 same as the second (and might actually have the same implementation).

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

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

942 **Symbolic Link**

943 Many implementations associate no attributes, including ownership with symbolic links.
944 Security experts encouraged consideration for defining these attributes as optional.
945 Consideration was given to changing *utime()* to allow modification of the times for a symbolic
946 link, or as an alternative adding an *lutime()* interface. Modifications to *chown()* were also
947 considered: allow changing symbolic link ownership or alternatively adding *lchown()*. As a
948 result of the problems encountered in defining attributes for symbolic links (and interfaces to
949 access/modify those attributes) and since implementations exist that do not associate these
950 attributes with symbolic links, only the file type bits in the *st_mode* member and the *st_size*
951 member of the **stat** structure are required to be applicable to symbolic links.

952 Historical implementations were followed when determining which interfaces should apply to
953 symbolic links. Interfaces that historically followed symbolic links include *chmod()*, *link()*, and
954 *utime()*. Interfaces that historically do not follow symbolic links include *chown()*, *lstat()*,
955 *readlink()*, *rename()*, *remove()*, *rmdir()*, and *unlink()*. IEEE Std 1003.1-200x deviates from
956 historical practice only in the case of *chown()*. Because there is no requirement that there be an
957 association of ownership with symbolic links, there was no point in requiring an interface to
958 change ownership. In addition, other implementations of symbolic links have modified *chown()*
959 to follow symbolic links.

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

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

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

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

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

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

- 992 1. Symbolic links specified to system calls that take filename arguments
- 993 2. Symbolic links specified as command line filename arguments to utilities that are not
994 performing a traversal of a file hierarchy
- 995 3. Symbolic links referencing files not of type directory, specified to utilities that are
996 performing a traversal of a file hierarchy
- 997 4. Symbolic links referencing files of type directory, specified to utilities that are performing a
998 traversal of a file hierarchy

999 *First Domain*

1000 The first domain is considered in earlier rationale.

1001 *Second Domain*

1002 The reason this category is restricted to utilities that are not traversing the file hierarchy is that
1003 some standard utilities take an option that specifies a hierarchical traversal, but by default
1004 operate on the arguments themselves. Generally, users specifying the option for a file hierarchy
1005 traversal wish to operate on a single, physical hierarchy, and therefore symbolic links, which
1006 may reference files outside of the hierarchy, are ignored. For example, *chown owner file* is a
1007 different operation from the same command with the **-R** option specified. In this example, the
1008 behavior of the command *chown owner file* is described here, while the behavior of the command
1009 *chown -R owner file* is described in the third and fourth domains.

1010 The general rule is that the utilities in this category follow symbolic links named as arguments.

1011 Exceptions in the second domain are:

- 1012 • The *mv* and *rm* utilities do not follow symbolic links named as arguments, but respectively
1013 attempt to rename or delete them.
- 1014 • The *ls* utility is also an exception to this rule. For compatibility with historical systems, when
1015 the **-R** option is not specified, the *ls* utility follows symbolic links named as arguments if the
1016 **-L** option is specified or if the **-F**, **-d**, or **-l** options are not specified. (If the **-L** option is
1017 specified, *ls* always follows symbolic links; it is the only utility where the **-L** option affects its
1018 behavior even though a tree walk is not being performed.)

1019 All other standard utilities, when not traversing a file hierarchy, always follow symbolic links
1020 named as arguments.

1021 Historical practice is that the **-h** option is specified if standard utilities are to act upon symbolic
1022 links instead of upon their targets. Examples of commands that have historically had a **-h** option
1023 for this purpose are the *chgrp*, *chown*, *file*, and *test* utilities.

1024 *Third Domain*

1025 The third domain is symbolic links, referencing files not of type directory, specified to utilities
1026 that are performing a traversal of a file hierarchy. (This includes symbolic links specified as
1027 command line filename arguments or encountered during the traversal.)

1028 The intention of the Shell and Utilities volume of IEEE Std 1003.1-200x is that the operation that
1029 the utility is performing is applied to the symbolic link itself, if that operation is applicable to
1030 symbolic links. The reason that the operation is not required is that symbolic links in some
1031 implementations do not have such attributes as a file owner, and therefore the *chown* operation
1032 would be meaningless. If symbolic links on the system have an owner, it is the intention that the
1033 utility *chown* cause the owner of the symbolic link to change. If symbolic links do not have an
1034 owner, the symbolic link should be ignored. Specifically, by default, no change should be made
1035 to the file referenced by the symbolic link.

1036 *Fourth Domain*

1037 The fourth domain is symbolic links referencing files of type directory, specified to utilities that
1038 are performing a traversal of a file hierarchy. (This includes symbolic links specified as
1039 command line filename arguments or encountered during the traversal.)

1040 Most standard utilities do not, by default, indirect into the file hierarchy referenced by the
1041 symbolic link. (The Shell and Utilities volume of IEEE Std 1003.1-200x uses the informal term
1042 *physical walk* to describe this case. The case where the utility does indirect through the symbolic
1043 link is termed a *logical walk*.)

1044 There are three reasons for the default to a physical walk:

- 1045 1. With very few exceptions, a physical walk has been the historical default on UNIX systems
1046 supporting symbolic links. Because some utilities (that is, *rm*) must default to a physical
1047 walk, regardless, changing historical practice in this regard would be confusing to users
1048 and needlessly incompatible.
- 1049 2. For systems where symbolic links have the historical file attributes (that is, *owner*, *group*,
1050 *mode*), defaulting to a logical traversal would require the addition of a new option to the
1051 commands to modify the attributes of the link itself. This is painful and more complex
1052 than the alternatives.
- 1053 3. There is a security issue with defaulting to a logical walk. Historically, the command
1054 *chown -R user file* has been safe for the superuser because *setuid* and *setgid* bits were lost
1055 when the ownership of the file was changed. If the walk were logical, changing ownership
1056 would no longer be safe because a user might have inserted a symbolic link pointing to any
1057 file in the tree. Again, this would necessitate the addition of an option to the commands
1058 doing hierarchy traversal to not indirect through the symbolic links, and historical scripts
1059 doing recursive walks would instantly become security problems. While this is mostly an
1060 issue for system administrators, it is preferable to not have different defaults for different
1061 classes of users.

1062 As consistently as possible, users may cause standard utilities performing a file hierarchy
1063 traversal to follow any symbolic links named on the command line, regardless of the type of file
1064 they reference, by specifying the *-H* (for half logical) option. This option is intended to make the
1065 command line name space look like the logical name space.

1066 As consistently as possible, users may cause standard utilities performing a file hierarchy
1067 traversal to follow any symbolic links named on the command line as well as any symbolic links
1068 encountered during the traversal, regardless of the type of file they reference, by specifying the
1069 *-L* (for logical) option. This option is intended to make the entire name space look like the
1070 logical name space.

1071 For consistency, implementors are encouraged to use the **-P** (for physical) flag to specify the
1072 physical walk in utilities that do logical walks by default for whatever reason. The only standard
1073 utilities that require the **-P** option are *cd* and *pwd*; see the note below.

1074 When one or more of the **-H**, **-L**, and **-P** flags can be specified, the last one specified determines
1075 the behavior of the utility. This permits users to alias commands so that the default behavior is a
1076 logical walk and then override that behavior on the command line.

1077 *Exceptions in the Third and Fourth Domains*

1078 The *ls* and *rm* utilities are exceptions to these rules. The *rm* utility never follows symbolic links
1079 and does not support the **-H**, **-L**, or **-P** options. Some historical versions of *ls* always followed
1080 symbolic links given on the command line whether the **-L** option was specified or not. Historical
1081 versions of *ls* did not support the **-H** option. In IEEE Std 1003.1-200x, unless one of the **-H** or **-L**
1082 options is specified, the *ls* utility only follows symbolic links to directories that are given as
1083 operands. The *ls* utility does not support the **-P** option.

1084 The Shell and Utilities volume of IEEE Std 1003.1-200x requires that the standard utilities *ls*, *find*,
1085 and *pax* detect infinite loops when doing logical walks; that is, a directory, or more commonly a
1086 symbolic link, that refers to an ancestor in the current file hierarchy. If the file system itself is
1087 corrupted, causing the infinite loop, it may be impossible to recover. Because *find* and *ls* are often
1088 used in system administration and security applications, they should attempt to recover and
1089 continue as best as they can. The *pax* utility should terminate because the archive it was creating
1090 is by definition corrupted. Other, less vital, utilities should probably simply terminate as well.
1091 Implementations are strongly encouraged to detect infinite loops in all utilities.

1092 Historical practice is shown in Table A-1 (on page 3320). The heading **SVID3** stands for the
1093 Third Edition of the System V Interface Definition.

1094 Historically, several shells have had built-in versions of the *pwd* utility. In some of these shells,
1095 *pwd* reported the physical path, and in others, the logical path. Implementations of the shell
1096 corresponding to IEEE Std 1003.1-200x must report the logical path by default. Earlier versions
1097 of IEEE Std 1003.1-200x did not require the *pwd* utility to be a built-in utility. Now that *pwd* is
1098 required to set an environment variable in the current shell execution environment, it must be a
1099 built-in utility.

1100 The *cd* command is required, by default, to treat the filename dot-dot logically. Implementors are
1101 required to support the **-P** flag in *cd* so that users can have their current environment handled
1102 physically. In 4.3 BSD, *chgrp* during tree traversal changed the group of the symbolic link, not
1103 the target. Symbolic links in 4.4 BSD do not have *owner*, *group*, *mode*, or other standard UNIX
1104 system file attributes.

1105

Table A-1 Historical Practice for Symbolic Links

Utility	SVID3	4.3 BSD	4.4 BSD	POSIX	Comments
1106				-L	Treat ". ." logically.
1107				-P	". ." physically.
1108					
1109			-H	-H	Follow command line symlinks.
1110			-h	-L	Follow symlinks.
1111	-h			-h	Affect the symlink.
1112					Affect the symlink.
1113			-H		Follow command line symlinks.
1114			-h		Follow symlinks.
1115			-H	-H	Follow command line symlinks.
1116			-h	-L	Follow symlinks.
1117	-h			-h	Affect the symlink.
1118			-H	-H	Follow command line symlinks.
1119			-h	-L	Follow symlinks.
1120	-L		-L		Follow symlinks.
1121			-H	-H	Follow command line symlinks.
1122			-h	-L	Follow symlinks.
1123	-h			-h	Affect the symlink.
1124			-H	-H	Follow command line symlinks.
1125			-h	-L	Follow symlinks.
1126	-follow		-follow		Follow symlinks.
1127	-s	-s	-s	-s	Create a symbolic link.
1128	-L	-L	-L	-L	Follow symlinks.
1129				-H	Follow command line symlinks.
1130					Operates on the symlink.
1131			-H	-H	Follow command line symlinks.
1132			-h	-L	Follow symlinks.
1133				-L	Printed path may contain symlinks.
1134				-P	Printed path will not contain symlinks.
1135					Operates on the symlink.
1136			-H		Follow command line symlinks.
1137		-h	-h		Follow symlinks.
1138	-h		-h	-h	Affect the symlink.

1139 **Synchronously-Generated Signal**

1140 Those signals that may be generated synchronously include SIGABRT, SIGBUS, SIGILL, SIGFPE,
1141 SIGPIPE, and SIGSEGV.

1142 Any signal sent via the *raise()* function or a *kill()* function targeting the current process is also
1143 considered synchronous.

1144 **System Call***

1145 The distinction between a *system call* and a *library routine* is an implementation detail that may
1146 differ between implementations and has thus been excluded from POSIX.1.

1147 See "Interface, Not Implementation" in the Preface.

1148 System Reboot

1149 A *system reboot* is an event initiated by an unspecified circumstance that causes all processes
1150 (other than special system processes) to be terminated in an implementation-defined manner,
1151 after which any changes to the state and contents of files created or written to by a Conforming
1152 POSIX.1 Application prior to the event are implementation-defined.

1153 Synchronized I/O Data (and File) Integrity Completion

1154 These terms specify that for synchronized read operations, pending writes must be successfully
1155 completed before the read operation can complete. This is motivated by two circumstances.
1156 Firstly, when synchronizing processes can access the same file, but not share common buffers
1157 (such as for a remote file system), this requirement permits the reading process to guarantee that
1158 it can read data written remotely. Secondly, having data written synchronously is insufficient to
1159 guarantee the order with respect to a subsequent write by a reading process, and thus this extra
1160 read semantic is necessary.

1161 Text File

1162 The term *text file* does not prevent the inclusion of control or other non-printable characters
1163 (other than NUL). Therefore, standard utilities that list text files as inputs or outputs are either
1164 able to process the special characters or they explicitly describe their limitations within their
1165 individual descriptions. The definition of *text file* has caused controversy. The only difference
1166 between text and binary files is that text files have lines of less than {LINE_MAX} bytes, with no
1167 NUL characters, each terminated by a <newline>. The definition allows a file with a single
1168 <newline>, but not a totally empty file, to be called a text file. If a file ends with an incomplete
1169 line it is not strictly a text file by this definition. The <newline> referred to in
1170 IEEE Std 1003.1-200x is not some generic line separator, but a single character; files created on
1171 systems where they use multiple characters for ends of lines are not portable to all conforming
1172 systems without some translation process unspecified by IEEE Std 1003.1-200x.

1173 Thread

1174 IEEE Std 1003.1-200x defines a thread to be a flow of control within a process. Each thread has a
1175 minimal amount of private state; most of the state associated with a process is shared among all
1176 of the threads in the process. While most multi-thread extensions to POSIX have taken this
1177 approach, others have made different decisions.

1178 **Note:** The choice to put threads within a process does not constrain implementations to implement
1179 threads in that manner. However, all functions have to behave as though threads share the
1180 indicated state information with the process from which they were created.

1181 Threads need to share resources in order to cooperate. Memory has to be widely shared between
1182 threads in order for the threads to cooperate at a fine level of granularity. Threads keep data
1183 structures and the locks protecting those data structures in shared memory. For a data structure
1184 to be usefully shared between threads, such structures should not refer to any data that can only
1185 be interpreted meaningfully by a single thread. Thus, any system resources that might be
1186 referred to in data structures need to be shared between all threads. File descriptors, pathnames,
1187 and pointers to stack variables are all things that programmers want to share between their
1188 threads. Thus, the file descriptor table, the root directory, the current working directory, and the
1189 address space have to be shared.

1190 Library implementations are possible as long as the effective behavior is as if system services
1191 invoked by one thread do not suspend other threads. This may be difficult for some library
1192 implementations on systems that do not provide asynchronous facilities.

1193 See Section B.2.9 (on page 3439) for additional rationale.

1194 **Thread ID**

1195 See Section B.2.9.2 (on page 3455) for additional rationale.

1196 **Thread-Safe Function**

1197 All functions required by IEEE Std 1003.1-200x need to be thread-safe; see Section A.4.16 (on
1198 page 3330) and Section B.2.9.1 (on page 3452) for additional rationale.

1199 **User Database**

1200 There are no references in IEEE Std 1003.1-200x to a *passwd file* or a *group file*, and there is no
1201 requirement that the *group* or *passwd* databases be kept in files containing editable text. Many
1202 large timesharing systems use *passwd* databases that are hashed for speed. Certain security
1203 classifications prohibit certain information in the *passwd* database from being publicly readable.

1204 The term *encoded* is used instead of *encrypted* in order to avoid the implementation connotations
1205 (such as reversibility or use of a particular algorithm) of the latter term.

1206 The *getgrent()*, *setgrent()*, *endgrent()*, *getpwent()*, *setpwent()*, and *endpwent()* functions are not
1207 included as part of the base standard because they provide a linear database search capability
1208 that is not generally useful (the *getpwuid()*, *getpwnam()*, *getgrgid()*, and *getgrnam()* functions are
1209 provided for keyed lookup) and because in certain distributed systems, especially those with
1210 different authentication domains, it may not be possible or desirable to provide an application
1211 with the ability to browse the system databases indiscriminately. They are provided on XSI-
1212 conformant systems due to their historical usage by many existing applications.

1213 A change from historical implementations is that the structures used by these functions have
1214 fields of the types **gid_t** and **uid_t**, which are required to be defined in the **<sys/types.h>** header.
1215 IEEE Std 1003.1-200x requires implementations to ensure that these types are defined by
1216 inclusion of **<grp.h>** and **<pwd.h>**, respectively, without imposing any name space pollution or
1217 errors from redefinition of types.

1218 IEEE Std 1003.1-200x is silent about the content of the strings containing user or group names.
1219 These could be digit strings. IEEE Std 1003.1-200x is also silent as to whether such digit strings
1220 bear any relationship to the corresponding (numeric) user or group ID.

1221 *Database Access*

1222 The thread-safe versions of the user and group database access functions return values in user-
1223 supplied buffers instead of possibly using static data areas that may be overwritten by each call.

1224 **Virtual Processor***

1225 The term *virtual processor* was chosen as a neutral term describing all kernel-level schedulable
1226 entities, such as processes, Mach tasks, or lightweight processes. Implementing threads using
1227 multiple processes as virtual processors, or implementing multiplexed threads above a virtual
1228 processor layer, should be possible, provided some mechanism has also been implemented for
1229 sharing state between processes or virtual processors. Many systems may also wish to provide
1230 implementations of threads on systems providing “shared processes” or “variable-weight
1231 processes”. It was felt that exposing such implementation details would severely limit the type
1232 of systems upon which the threads interface could be supported and prevent certain types of
1233 valid implementations. It was also determined that a virtual processor interface was out of the
1234 scope of the Rationale (Informative) volume of IEEE Std 1003.1-200x.

1235 **XSI**

1236 This is introduced to allow IEEE Std 1003.1-200x to be adopted as an IEEE standard and an Open
1237 Group Technical Standard, serving both the POSIX and the Single UNIX Specification in a core
1238 set of volumes.

1239 The term *XSI* has been used for 10 years in connection with the XPG series and the first and
1240 second versions of the base volumes of the Single UNIX Specification. The XSI margin code was
1241 introduced to denote the extended or more restrictive semantics beyond POSIX that are
1242 applicable to UNIX systems.

1243 **A.4 General Concepts**1244 **A.4.1 Concurrent Execution**

1245 There is no additional rationale provided for this section.

1246 **A.4.2 Directory Protection**

1247 There is no additional rationale provided for this section.

1248 **A.4.3 Extended Security Controls**

1249 Allowing an implementation to define extended security controls enables the use of
1250 IEEE Std 1003.1-200x in environments that require different or more rigorous security than that
1251 provided in POSIX.1. Extensions are allowed in two areas: privilege and file access permissions.
1252 The semantics of these areas have been defined to permit extensions with reasonable, but not
1253 exact, compatibility with all existing practices. For example, the elimination of the superuser
1254 definition precludes identifying a process as privileged or not by virtue of its effective user ID.

1255 **A.4.4 File Access Permissions**

1256 A process should not try to anticipate the result of an attempt to access data by *a priori* use of
1257 these rules. Rather, it should make the attempt to access data and examine the return value (and
1258 possibly *errno* as well), or use *access()*. An implementation may include other security
1259 mechanisms in addition to those specified in POSIX.1, and an access attempt may fail because of
1260 those additional mechanisms, even though it would succeed according to the rules given in this
1261 section. (For example, the user's security level might be lower than that of the object of the access
1262 attempt.) The supplementary group IDs provide another reason for a process to not attempt to
1263 anticipate the result of an access attempt.

1264 **A.4.5 File Hierarchy**

1265 Though the file hierarchy is commonly regarded to be a tree, POSIX.1 does not define it as such
1266 for three reasons:

- 1267 1. Links may join branches.
- 1268 2. In some network implementations, there may be no single absolute root directory; see
1269 *pathname resolution*.
- 1270 3. With symbolic links, the file system need not be a tree or even a directed acyclic graph.

1271 **A.4.6 Filenames**

1272 Historically, certain filenames have been reserved. This list includes **core**, **/etc/passwd**, and so on. Conforming applications should avoid these.

1274 Most historical implementations prohibit case folding in filenames; that is, treating uppercase and lowercase alphabetic characters as identical. However, some consider case folding desirable:

- 1276 • For user convenience
- 1277 • For ease-of-implementation of the POSIX.1 interface as a hosted system on some popular
- 1278 operating systems

1279 Variants, such as maintaining case distinctions in filenames, but ignoring them in comparisons, have been suggested. Methods of allowing escaped characters of the case opposite the default have been proposed.

1282 Many reasons have been expressed for not allowing case folding, including:

- 1283 • No solid evidence has been produced as to whether case-sensitivity or case-insensitivity is
- 1284 more convenient for users.
- 1285 • Making case-insensitivity a POSIX.1 implementation option would be worse than either
- 1286 having it or not having it, because:
 - 1287 — More confusion would be caused among users.
 - 1288 — Application developers would have to account for both cases in their code.
 - 1289 — POSIX.1 implementors would still have other problems with native file systems, such as
 - 1290 short or otherwise constrained filenames or pathnames, and the lack of hierarchical
 - 1291 directory structure.
- 1292 • Case folding is not easily defined in many European languages, both because many of them
- 1293 use characters outside the US ASCII alphabetic set, and because:
 - 1294 — In Spanish, the digraph "ll" is considered to be a single letter, the capitalized form of
 - 1295 which may be either "Ll" or "LL", depending on context.
 - 1296 — In French, the capitalized form of a letter with an accent may or may not retain the accent,
 - 1297 depending on the country in which it is written.
 - 1298 — In German, the sharp ess may be represented as a single character resembling a Greek
 - 1299 beta (β) in lowercase, but as the digraph "SS" in uppercase.
 - 1300 — In Greek, there are several lowercase forms of some letters; the one to use depends on its
 - 1301 position in the word. Arabic has similar rules.
- 1302 • Many East Asian languages, including Japanese, Chinese, and Korean, do not distinguish
- 1303 case and are sometimes encoded in character sets that use more than one byte per character.
- 1304 • Multiple character codes may be used on the same machine simultaneously. There are
- 1305 several ISO character sets for European alphabets. In Japan, several Japanese character codes
- 1306 are commonly used together, sometimes even in filenames; this is evidently also the case in
- 1307 China. To handle case insensitivity, the kernel would have to at least be able to distinguish
- 1308 for which character sets the concept made sense.
- 1309 • The file system implementation historically deals only with bytes, not with characters, except
- 1310 for slash and the null byte.
- 1311 • The purpose of POSIX.1 is to standardize the common, existing definition, not to change it.
- 1312 Mandating case-insensitivity would make all historical implementations non-standard.

1313 • Not only the interface, but also application programs would need to change, counter to the
1314 purpose of having minimal changes to existing application code.

1315 • At least one of the original developers of the UNIX system has expressed objection in the
1316 strongest terms to either requiring case-insensitivity or making it an option, mostly on the
1317 basis that POSIX.1 should not hinder portability of application programs across related
1318 implementations in order to allow compatibility with unrelated operating systems.

1319 Two proposals were entertained regarding case folding in filenames:

1320 1. Remove all wording that previously permitted case folding.

1321 Rationale Case folding is inconsistent with portable filename character set definition
1322 and filename definition (all characters except slash and null). No known
1323 implementations allowing all characters except slash and null also do case
1324 folding.

1325 2. Change “though this practice is not recommended:” to “although this practice is strongly
1326 discouraged.”

1327 Rationale If case folding must be included in POSIX.1, the wording should be stronger
1328 to discourage the practice.

1329 The consensus selected the first proposal. Otherwise, a conforming application would have to
1330 assume that case folding would occur when it was not wanted, but that it would not occur when
1331 it was wanted.

1332 **A.4.7 File Times Update**

1333 This section reflects the actions of historical implementations. The times are not updated
1334 immediately, but are only marked for update by the functions. An implementation may update
1335 these times immediately.

1336 The accuracy of the time update values is intentionally left unspecified so that systems can
1337 control the bandwidth of a possible covert channel.

1338 The wording was carefully chosen to make it clear that there is no requirement that the
1339 conformance document contain information that might incidentally affect file update times. Any
1340 function that performs pathname resolution might update several *st_atime* fields. Functions such
1341 as *getpwnam()* and *getgrnam()* might update the *st_atime* field of some specific file or files. It
1342 is intended that these are not required to be documented in the conformance document, but they
1343 should appear in the system documentation.

1344 **A.4.8 Host and Network Byte Order**

1345 There is no additional rationale provided for this section.

1346 **A.4.9 Measurement of Execution Time**

1347 The methods used to measure the execution time of processes and threads, and the precision of
1348 these measurements, may vary considerably depending on the software architecture of the
1349 implementation, and on the underlying hardware. Implementations can also make tradeoffs
1350 between the scheduling overhead and the precision of the execution time measurements.
1351 IEEE Std 1003.1-200x does not impose any requirement on the accuracy of the execution time; it
1352 instead specifies that the measurement mechanism and its precision are implementation-
1353 defined.

1354 **A.4.10 Memory Synchronization**

1355 In older multi-processors, access to memory by the processors was strictly multiplexed. This
 1356 meant that a processor executing program code interrogates or modifies memory in the order
 1357 specified by the code and that all the memory operation of all the processors in the system
 1358 appear to happen in some global order, though the operation histories of different processors are
 1359 interleaved arbitrarily. The memory operations of such machines are said to be sequentially
 1360 consistent. In this environment, threads can synchronize using ordinary memory operations. For
 1361 example, a producer thread and a consumer thread can synchronize access to a circular data
 1362 buffer as follows:

```

1363     int rdptr = 0;
1364     int wrptr = 0;
1365     data_t buf[BUFSIZE];

1366     Thread 1:
1367         while (work_to_do) {
1368             int next;

1369             buf[wrptr] = produce();
1370             next = (wrptr + 1) % BUFSIZE;
1371             while (rdptr == next)
1372                 ;
1373             wrptr = next;
1374         }

1375     Thread 2:
1376         while (work_to_do) {
1377             while (rdptr == wrptr)
1378                 ;
1379             consume(buf[rdptr]);
1380             rdptr = (rdptr + 1) % BUFSIZE;
1381         }

```

1382 In modern multi-processors, these conditions are relaxed to achieve greater performance. If one
 1383 processor stores values in location A and then location B, then other processors loading data
 1384 from location B and then location A may see the new value of B but the old value of A. The
 1385 memory operations of such machines are said to be weakly ordered. On these machines, the
 1386 circular buffer technique shown in the example will fail because the consumer may see the new
 1387 value of *wrptr* but the old value of the data in the buffer. In such machines, synchronization can
 1388 only be achieved through the use of special instructions that enforce an order on memory
 1389 operations. Most high-level language compilers only generate ordinary memory operations to
 1390 take advantage of the increased performance. They usually cannot determine when memory
 1391 operation order is important and generate the special ordering instructions. Instead, they rely on
 1392 the programmer to use synchronization primitives correctly to ensure that modifications to a
 1393 location in memory are ordered with respect to modifications and/or access to the same location
 1394 in other threads. Access to read-only data need not be synchronized. The resulting program is
 1395 said to be data race-free.

1396 Synchronization is still important even when accessing a single primitive variable (for example,
 1397 an integer). On machines where the integer may not be aligned to the bus data width or be larger
 1398 than the data width, a single memory load may require multiple memory cycles. This means
 1399 that it may be possible for some parts of the integer to have an old value while other parts have a
 1400 newer value. On some processor architectures this cannot happen, but portable programs cannot
 1401 rely on this.

1402 In summary, a portable multi-threaded program, or a multi-process program that shares
1403 writable memory between processes, has to use the synchronization primitives to synchronize
1404 data access. It cannot rely on modifications to memory being observed by other threads in the
1405 order written in the program or even on modification of a single variable being seen atomically.

1406 Conforming applications may only use the functions listed to synchronize threads of control
1407 with respect to memory access. There are many other candidates for functions that might also be
1408 used. Examples are: signal sending and reception, or pipe writing and reading. In general, any
1409 function that allows one thread of control to wait for an action caused by another thread of
1410 control is a candidate. IEEE Std 1003.1-200x does not require these additional functions to
1411 synchronize memory access since this would imply the following:

- 1412 • All these functions would have to be recognized by advanced compilation systems so that
1413 memory operations and calls to these functions are not reordered by optimization.
- 1414 • All these functions would potentially have to have memory synchronization instructions
1415 added, depending on the particular machine.
- 1416 • The additional functions complicate the model of how memory is synchronized and make
1417 automatic data race detection techniques impractical.

1418 Formal definitions of the memory model were rejected as unreadable by the vast majority of
1419 programmers. In addition, most of the formal work in the literature has concentrated on the
1420 memory as provided by the hardware as opposed to the application programmer through the
1421 compiler and runtime system. It was believed that a simple statement intuitive to most
1422 programmers would be most effective. IEEE Std 1003.1-200x defines functions that can be used
1423 to synchronize access to memory, but it leaves open exactly how one relates those functions to
1424 the semantics of each function as specified elsewhere in IEEE Std 1003.1-200x.
1425 IEEE Std 1003.1-200x also does not make a formal specification of the partial ordering in time
1426 that the functions can impose, as that is implied in the description of the semantics of each
1427 function. It simply states that the programmer has to ensure that modifications do not occur
1428 “simultaneously” with other access to a memory location.

1429 **A.4.11 Pathname Resolution**

1430 It is necessary to differentiate between the definition of pathname and the concept of pathname
1431 resolution with respect to the handling of trailing slashes. By specifying the behavior here, it is
1432 not possible to provide an implementation that is conforming but extends all interfaces that
1433 handle pathnames to also handle strings that are not legal pathnames (because they have trailing
1434 slashes).

1435 Pathnames that end with one or more trailing slash characters must refer to directory paths.
1436 Previous versions of IEEE Std 1003.1-200x were not specific about the distinction between
1437 trailing slashes on files and directories, and both were permitted.

1438 Two types of implementation have been prevalent; those that ignored trailing slash characters
1439 on all pathnames regardless, and those that only permitted them only on existing directories.

1440 IEEE Std 1003.1-200x requires that a pathname with a trailing slash character be treated as if it
1441 had a trailing " / . " everywhere.

1442 Note that this change does not break any conforming applications; since there were two
1443 different types of implementation, no application could have portably depended on either
1444 behavior. This change does however require some implementations to be altered to remain
1445 compliant. Substantial discussion over a three-year period has shown that the benefits to
1446 application developers outweighs the disadvantages for some vendors.

1447 On a historical note, some early applications automatically appended a '/' to every path.
 1448 Rather than fix the applications, the system implementation was modified to accept this
 1449 behavior by ignoring any trailing slash.

1450 Each directory has exactly one parent directory which is represented by the name **dot-dot** in the
 1451 first directory. No other directory, regardless of linkages established by symbolic links, is
 1452 considered the parent directory by IEEE Std 1003.1-200x.

1453 There are two general categories of interfaces involving pathname resolution: those that follow
 1454 the symbolic link, and those that do not. There are several exceptions to this rule; for example,
 1455 *open(path,O_CREAT|O_EXCL)* will fail when *path* names a symbolic link. However, in all other
 1456 situations, the *open()* function will follow the link.

1457 What the filename **dot-dot** refers to relative to the root directory is implementation-defined. In
 1458 Version 7 it refers to the root directory itself; this is the behavior mentioned in
 1459 IEEE Std 1003.1-200x. In some networked systems the construction *././hostname/* is used to refer
 1460 to the root directory of another host, and POSIX.1 permits this behavior.

1461 Other networked systems use the construct *//hostname* for the same purpose; that is, a double
 1462 initial slash is used. There is a potential problem with existing applications that create full
 1463 pathnames by taking a trunk and a relative pathname and making them into a single string
 1464 separated by '/', because they can accidentally create networked pathnames when the trunk is
 1465 '/'. This practice is not prohibited because such applications can be made to conform by
 1466 simply changing to use "//" as a separator instead of '/':

- 1467 • If the trunk is '/', the full pathname will begin with "///" (the initial '/' and the
 1468 separator "//"). This is the same as '/', which is what is desired. (This is the general case
 1469 of making a relative pathname into an absolute one by prefixing with "///" instead of '/'.)
- 1470 • If the trunk is "/A", the result is "/A//..."; since non-leading sequences of two or more
 1471 slashes are treated as a single slash, this is equivalent to the desired "/A/...".
- 1472 • If the trunk is "//A", the implementation-defined semantics will apply. (The multiple slash
 1473 rule would apply.)

1474 Application developers should avoid generating pathnames that start with "//".
 1475 Implementations are strongly encouraged to avoid using this special interpretation since a
 1476 number of applications currently do not follow this practice and may inadvertently generate
 1477 "///...".

1478 The term *root directory* is only defined in POSIX.1 relative to the process. In some
 1479 implementations, there may be no absolute root directory. The initialization of the root directory
 1480 of a process is implementation-defined.

1481 **A.4.12 Process ID Reuse**

1482 There is no additional rationale provided for this section.

1483 **A.4.13 Scheduling Policy**

1484 There is no additional rationale provided for this section.

1485 **A.4.14 Seconds Since the Epoch**

1486 Coordinated Universal Time (UTC) includes leap seconds. However, in POSIX time (seconds
1487 since the Epoch), leap seconds are ignored (not applied) to provide an easy and compatible
1488 method of computing time differences. Broken-down POSIX time is therefore not necessarily
1489 UTC, despite its appearance.

1490 As of September 2000, 24 leap seconds had been added to UTC since the Epoch, 1 January, 1970.
1491 Historically, one leap second is added every 15 months on average, so this offset can be expected
1492 to grow steadily with time.

1493 Most systems' notion of "time" is that of a continuously increasing value, so this value should
1494 increase even during leap seconds. However, not only do most systems not keep track of leap
1495 seconds, but most systems are probably not synchronized to any standard time reference.
1496 Therefore, it is inappropriate to require that a time represented as seconds since the Epoch
1497 precisely represent the number of seconds between the referenced time and the Epoch.

1498 It is sufficient to require that applications be allowed to treat this time as if it represented the
1499 number of seconds between the referenced time and the Epoch. It is the responsibility of the
1500 vendor of the system, and the administrator of the system, to ensure that this value represents
1501 the number of seconds between the referenced time and the Epoch as closely as necessary for the
1502 application being run on that system.

1503 It is important that the interpretation of time names and *seconds since the Epoch* values be
1504 consistent across conforming systems; that is, it is important that all conforming systems
1505 interpret "536 457 599 seconds since the Epoch" as 59 seconds, 59 minutes, 23 hours 31 December
1506 1986, regardless of the accuracy of the system's idea of the current time. The expression is given
1507 to assure a consistent interpretation, not to attempt to specify the calendar. The relationship
1508 between *tm_yday* and the day of week, day of month, and month is presumed to be specified
1509 elsewhere and is not given in POSIX.1.

1510 Consistent interpretation of *seconds since the Epoch* can be critical to certain types of distributed
1511 applications that rely on such timestamps to synchronize events. The accrual of leap seconds in
1512 a time standard is not predictable. The number of leap seconds since the Epoch will likely
1513 increase. POSIX.1 is more concerned about the synchronization of time between applications of
1514 astronomically short duration.

1515 Note that *tm_yday* is zero-based, not one-based, so the day number in the example above is 364.
1516 Note also that the division is an integer division (discarding remainder) as in the C language.

1517 Note also that the meaning of *gmtime()*, *localtime()*, and *mktime()* is specified in terms of this
1518 expression. However, the ISO C standard computes *tm_yday* from *tm_mday*, *tm_mon*, and
1519 *tm_year* in *mktime()*. Because it is stated as a (bidirectional) relationship, not a function, and
1520 because the conversion between month-day-year and day-of-year dates is presumed well known
1521 and is also a relationship, this is not a problem.

1522 Implementations that implement **time_t** as a signed 32-bit integer will overflow in 2 038. The
1523 data size for **time_t** is as per the ISO C standard definition, which is implementation-defined.

1524 See also **Epoch** (on page 3306).

1525 The topic of whether seconds since the Epoch should account for leap seconds has been debated
1526 on a number of occasions, and each time consensus was reached (with acknowledged dissent
1527 each time) that the majority of users are best served by treating all days identically. (That is, the

1528 majority of applications were judged to assume a single length—as measured in seconds since
 1529 the Epoch—for all days. Thus, leap seconds are not applied to seconds since the Epoch.) Those
 1530 applications which do care about leap seconds can determine how to handle them in whatever
 1531 way those applications feel is best. This was particularly emphasized because there was
 1532 disagreement about what the best way of handling leap seconds might be. It is a practical
 1533 impossibility to mandate that a conforming implementation must have a fixed relationship to
 1534 any particular official clock (consider isolated systems, or systems performing “reruns” by
 1535 setting the clock to some arbitrary time).

1536 Note that as a practical consequence of this, the length of a second as measured by some external
 1537 standard is not specified. This unspecified second is nominally equal to an International System
 1538 (SI) second in duration. Applications must be matched to a system that provides the particular
 1539 handling of external time in the way required by the application.

1540 **A.4.15 Semaphore**

1541 There is no additional rationale provided for this section.

1542 **A.4.16 Thread-Safety**

1543 Where the interface of a function required by IEEE Std 1003.1-200x precludes thread-safety, an
 1544 alternate form that shall be thread-safe is provided. The names of these thread-safe forms are the
 1545 same as the non-thread-safe forms with the addition of the suffix “_r”. The suffix “_r” is
 1546 historical, where the ‘r’ stood for “reentrant”.

1547 In some cases, thread-safety is provided by restricting the arguments to an existing function.

1548 See also Section B.2.9.1 (on page 3452).

1549 **A.4.17 Tracing**

1550 Refer to Section B.2.11 (on page 3468).

1551 **A.4.18 Treatment of Error Conditions for Mathematical Functions**

1552 There is no additional rationale provided for this section. |

1553 **A.4.19 Treatment of NaN Arguments for Mathematical Functions**

1554 There is no additional rationale provided for this section. |

1555 **A.4.20 Utility**

1556 There is no additional rationale provided for this section.

1557 **A.4.21 Variable Assignment**

1558 There is no additional rationale provided for this section.

1559 A.5 File Format Notation

1560 The notation for spaces allows some flexibility for application output. Note that an empty
 1561 character position in *format* represents one or more <blank>s on the output (not *white space*,
 1562 which can include <newline>s). Therefore, another utility that reads that output as its input
 1563 must be prepared to parse the data using *scanf()*, *awk*, and so on. The 'Δ' character is used when
 1564 exactly one <space> is output.

1565 The treatment of integers and spaces is different from the *printf()* function in that they can be
 1566 surrounded with <blank>s. This was done so that, given a format such as:

```
1567     "%d\n", <foo>
```

1568 the implementation could use a *printf()* call such as:

```
1569     printf("%6d\n", foo);
```

1570 and still conform. This notation is thus somewhat like *scanf()* in addition to *printf()*.

1571 The *printf()* function was chosen as a model because most of the standard developers were
 1572 familiar with it. One difference from the C function *printf()* is that the l and h conversion
 1573 specifier characters are not used. As expressed by the Shell and Utilities volume of
 1574 IEEE Std 1003.1-200x, there is no differentiation between decimal values for type **int**, type **long**,
 1575 or type **short**. The conversion specifications %d or %i should be interpreted as an arbitrary
 1576 length sequence of digits. Also, no distinction is made between single precision and double
 1577 precision numbers (**float** or **double** in C). These are simply referred to as floating-point numbers.

1578 Many of the output descriptions in the Shell and Utilities volume of IEEE Std 1003.1-200x use the
 1579 term *line*, such as:

```
1580     "%s", <input line>
```

1581 Since the definition of *line* includes the trailing <newline> already, there is no need to include a
 1582 '\n' in the format; a double <newline> would otherwise result.

1583 A.6 Character Set

1584 A.6.1 Portable Character Set

1585 The portable character set is listed in full so there is no dependency on the ISO/IEC 646:1991
 1586 standard (or historically ASCII) encoded character set, although the set is identical to the
 1587 characters defined in the International Reference version of the ISO/IEC 646:1991 standard.

1588 IEEE Std 1003.1-200x poses no requirement that multiple character sets or codesets be
 1589 supported, leaving this as a marketing differentiation for implementors. Although multiple
 1590 charmap files are supported, it is the responsibility of the implementation to provide the file(s);
 1591 if only one is provided, only that one will be accessible using the *localedef -f* option.

1592 The statement about invariance in codesets for the portable character set is worded to avoid
 1593 precluding implementations where multiple incompatible codesets are available (for instance,
 1594 ASCII and EBCDIC). The standard utilities cannot be expected to produce predictable results if
 1595 they access portable characters that vary on the same implementation.

1596 Not all character sets need include the portable character set, but each locale must include it. For
 1597 example, a Japanese-based locale might be supported by a mixture of character sets: JIS X 0201
 1598 Roman (a Japanese version of the ISO/IEC 646:1991 standard), JIS X 0208, and JIS X 0201
 1599 Katakana. Not all of these character sets include the portable characters, but at least one does
 1600 (JIS X 0201 Roman).

1601 **A.6.2 Character Encoding**

1602 Encoding mechanisms based on single shifts, such as the EUC encoding used in some Asian and
 1603 other countries, can be supported via the current charmap mechanism. With single-shift
 1604 encoding, each character is preceded by a shift code (SS2 or SS3). A complete EUC code,
 1605 consisting of the portable character set (G0) and up to three additional character sets (G1, G2,
 1606 G3), can be described using the current charmap mechanism; the encoding for each character in
 1607 additional character sets G2 and G3 must then include their single-shift code. Other mechanisms
 1608 to support locales based on encoding mechanisms such as locking shift are not addressed by this
 1609 volume of IEEE Std 1003.1-200x.

1610 **A.6.3 C Language Wide-Character Codes**

1611 There is no additional rationale provided for this section.

1612 **A.6.4 Character Set Description File**

1613 IEEE PASC Interpretation 1003.2 #196 is applied, removing three lines of text dealing with |
 1614 ranges of symbolic names using position constant values which had been erroneously included |
 1615 in the final 1003.2b draft. |

1616 *A.6.4.1 State-Dependent Character Encodings*

1617 A requirement was considered that would force utilities to eliminate any redundant locking
 1618 shifts, but this was left as a quality of implementation issue.

1619 This change satisfies the following requirement from the ISO POSIX-2:1993 standard, Annex
 1620 H.1:

1621 *The support of state-dependent (shift encoding) character sets should be addressed fully. See*
 1622 *descriptions of these in the Base Definitions volume of IEEE Std 1003.1-200x, Section 6.2, Character*
 1623 *Encoding. If such character encodings are supported, it is expected that this will impact the Base*
 1624 *Definitions volume of IEEE Std 1003.1-200x, Section 6.2, Character Encoding, the Base Definitions*
 1625 *volume of IEEE Std 1003.1-200x, Chapter 7, Locale, the Base Definitions volume of*
 1626 *IEEE Std 1003.1-200x, Chapter 9, Regular Expressions , and the comm, cut, diff, grep, head, join,*
 1627 *paste, and tail utilities.*

1628 The character set description file provides:

- 1629 • The capability to describe character set attributes (such as collation order or character
 1630 classes) independent of character set encoding, and using only the characters in the portable
 1631 character set. This makes it possible to create generic *localedef* source files for all codesets that
 1632 share the portable character set (such as the ISO 8859 family or IBM Extended ASCII).
- 1633 • Standardized symbolic names for all characters in the portable character set, making it
 1634 possible to refer to any such character regardless of encoding.

1635 Implementations are free to choose their own symbolic names, as long as the names identified
 1636 by this volume of IEEE Std 1003.1-200x are also defined; this provides support for already
 1637 existing “character names”.

1638 The names selected for the members of the portable character set follow the
 1639 ISO/IEC 8859-1:1998 standard and the ISO/IEC 10646-1:2000 standard. However, several
 1640 commonly used UNIX system names occur as synonyms in the list:

- 1641 • The historical UNIX system names are used for control characters.

- 1642 • The word “slash” is given in addition to “solidus”.
 - 1643 • The word “backslash” is given in addition to “reverse-solidus”.
 - 1644 • The word “hyphen” is given in addition to “hyphen-minus”.
 - 1645 • The word “period” is given in addition to “full-stop”.
 - 1646 • For digits, the word “digit” is eliminated.
 - 1647 • For letters, the words “Latin Capital Letter” and “Latin Small Letter” are eliminated.
 - 1648 • The words “left brace” and “right brace” are given in addition to “left-curly-bracket” and
 - 1649 “right-curly-bracket”.
 - 1650 • The names of the digits are preferred over the numbers to avoid possible confusion between
 - 1651 ‘0’ and ‘o’, and between ‘1’ and ‘l’ (one and the letter ell).
- 1652 The names for the control characters in the Base Definitions volume of IEEE Std 1003.1-200x,
- 1653 Chapter 6, Character Set were taken from the ISO/IEC 4873: 1991 standard.
- 1654 The charmap file was introduced to resolve problems with the portability of, especially, *localedef*
- 1655 sources. IEEE Std 1003.1-200x assumes that the portable character set is constant across all
- 1656 locales, but does not prohibit implementations from supporting two incompatible codings, such
- 1657 as both ASCII and EBCDIC. Such dual-support implementations should have all charmaps and
- 1658 *localedef* sources encoded using one portable character set, in effect cross-compiling for the other
- 1659 environment. Naturally, charmaps (and *localedef* sources) are only portable without
- 1660 transformation between systems using the same encodings for the portable character set. They
- 1661 can, however, be transformed between two sets using only a subset of the actual characters (the
- 1662 portable character set). However, the particular coded character set used for an application or an
- 1663 implementation does not necessarily imply different characteristics or collation; on the contrary,
- 1664 these attributes should in many cases be identical, regardless of codeset. The charmap provides
- 1665 the capability to define a common locale definition for multiple codesets (the same *localedef*
- 1666 source can be used for codesets with different extended characters; the ability in the charmap to
- 1667 define empty names allows for characters missing in certain codesets).
- 1668 The `<escape_char>` declaration was added at the request of the international community to ease
- 1669 the creation of portable charmap files on terminals not implementing the default backslash
- 1670 escape. The `<comment_char>` declaration was added at the request of the international
- 1671 community to eliminate the potential confusion between the number sign and the pound sign.
- 1672 The octal number notation with no leading zero required was selected to match those of *awk* and
- 1673 *tr* and is consistent with that used by *localedef*. To avoid confusion between an octal constant
- 1674 and the back-references used in *localedef* source, the octal, hexadecimal, and decimal constants
- 1675 shall contain at least two digits. As single-digit constants are relatively rare, this should not
- 1676 impose any significant hardship. Provision is made for more digits to account for systems in
- 1677 which the byte size is larger than 8 bits. For example, a Unicode (ISO/IEC 10646-1:2000
- 1678 standard) system that has defined 16-bit bytes may require six octal, four hexadecimal, and five
- 1679 decimal digits.
- 1680 The decimal notation is supported because some newer international standards define character
- 1681 values in decimal, rather than in the old column/row notation.
- 1682 The charmap identifies the coded character sets supported by an implementation. At least one
- 1683 charmap shall be provided, but no implementation is required to provide more than one.
- 1684 Likewise, implementations can allow users to generate new charmaps (for instance, for a new
- 1685 version of the ISO 8859 family of coded character sets), but does not have to do so. If users are
- 1686 allowed to create new charmaps, the system documentation describes the rules that apply (for
- 1687 instance, “only coded character sets that are supersets of the ISO/IEC 646: 1991 standard IRV, no

1688 multi-byte characters’’).

1689 This addition of the **WIDTH** specification satisfies the following requirement from the
1690 ISO POSIX-2:1993 standard, Annex H.1:

1691 (9) *The definition of column position relies on the implementation’s knowledge of the integral width*
1692 *of the characters. The charmap or LC_CTYPE locale definitions should be enhanced to allow*
1693 *application specification of these widths.*

1694 The character “width” information was first considered for inclusion under *LC_CTYPE* but was
1695 moved because it is more closely associated with the information in the *charmap* than
1696 information in the locale source (cultural conventions information). Concerns were raised that
1697 formalizing this type of information is moving the locale source definition from the codeset-
1698 independent entity that it was designed to be to a repository of codeset-specific information. A
1699 similar issue occurred with the `<code_set_name>`, `<mb_cur_max>`, and `<mb_cur_min>`
1700 information, which was resolved to reside in the *charmap* definition.

1701 The width definition was added to the IEEE P1003.2b draft standard with the intent that the
1702 *wcswidth()* and/or *wcwidth()* functions (currently specified in the System Interfaces volume of
1703 IEEE Std 1003.1-200x) be the mechanism to retrieve the character width information.

1704 **A.7 Locale**

1705 **A.7.1 General**

1706 The description of locales is based on work performed in the UniForum Technical Committee
1707 Subcommittee on Internationalization. Wherever appropriate, keywords are taken from the
1708 ISO C standard or the X/Open Portability Guide.

1709 The value used to specify a locale with environment variables is the name specified as the *name*
1710 operand to the *localedef* utility when the locale was created. This provides a verifiable method to
1711 create and invoke a locale.

1712 The “object” definitions need not be portable, as long as “source” definitions are. Strictly
1713 speaking, source definitions are portable only between implementations using the same
1714 character set(s). Such source definitions, if they use symbolic names only, easily can be ported
1715 between systems using different codesets, as long as the characters in the portable character set
1716 (see the Base Definitions volume of IEEE Std 1003.1-200x, Section 6.1, Portable Character Set)
1717 have common values between the codesets; this is frequently the case in historical
1718 implementations. Of source, this requires that the symbolic names used for characters outside
1719 the portable character set be identical between character sets. The definition of symbolic names
1720 for characters is outside the scope of IEEE Std 1003.1-200x, but is certainly within the scope of
1721 other standards organizations.

1722 Applications can select the desired locale by invoking the *setlocale()* function (or equivalent)
1723 with the appropriate value. If the function is invoked with an empty string, the value of the
1724 corresponding environment variable is used. If the environment variable is not set or is set to the
1725 empty string, the implementation sets the appropriate environment as defined in the Base
1726 Definitions volume of IEEE Std 1003.1-200x, Chapter 8, Environment Variables.

1727 A.7.2 POSIX Locale

1728 The POSIX locale is equal to the C locale. To avoid being classified as a C-language function, the
1729 name has been changed to the POSIX locale; the environment variable value can be either
1730 "POSIX" or, for historical reasons, "C".

1731 The POSIX definitions mirror the historical UNIX system behavior.

1732 The use of symbolic names for characters in the tables does not imply that the POSIX locale must
1733 be described using symbolic character names, but merely that it may be advantageous to do so.

1734 A.7.3 Locale Definition

1735 The decision to separate the file format from the *localedef* utility description was only partially
1736 editorial. Implementations may provide other interfaces than *localedef*. Requirements on “the
1737 utility”, mostly concerning error messages, are described in this way because they are meant to
1738 affect the other interfaces implementations may provide as well as *localedef*.

1739 The text about POSIX2_LOCALEDEF does not mean that internationalization is optional; only
1740 that the functionality of the *localedef* utility is. REs, for instance, must still be able to recognize,
1741 for example, character class expressions such as "[[:alpha:]]". A possible analogy is with
1742 an applications development environment; while all conforming implementations must be
1743 capable of executing applications, not all need to have the development environment installed.
1744 The assumption is that the capability to modify the behavior of utilities (and applications) via
1745 locale settings must be supported. If the *localedef* utility is not present, then the only choice is to
1746 select an existing (presumably implementation-documented) locale. An implementation could,
1747 for example, choose to support only the POSIX locale, which would in effect limit the amount of
1748 changes from historical implementations quite drastically. The *localedef* utility is still required,
1749 but would always terminate with an exit code indicating that no locale could be created.
1750 Supported locales must be documented using the syntax defined in this chapter. (This ensures
1751 that users can accurately determine what capabilities are provided. If the implementation
1752 decides to provide additional capabilities to the ones in this chapter, that is already provided
1753 for.)

1754 If the option is present (that is, locales can be created), then the *localedef* utility must be capable
1755 of creating locales based on the syntax and rules defined in this chapter. This does not mean that
1756 the implementation cannot also provide alternate means for creating locales.

1757 The octal, decimal, and hexadecimal notations are the same employed by the charmap facility
1758 (see the Base Definitions volume of IEEE Std 1003.1-200x, Section 6.4, Character Set Description
1759 File). To avoid confusion between an octal constant and a back-reference, the octal, hexadecimal,
1760 and decimal constants must contain at least two digits. As single-digit constants are relatively
1761 rare, this should not impose any significant hardship. Provision is made for more digits to
1762 account for systems in which the byte size is larger than 8 bits. For example, a Unicode (see the
1763 ISO/IEC 10646-1:2000 standard) system that has defined 16-bit bytes may require six octal, four
1764 hexadecimal, and five decimal digits. As with the charmap file, multi-byte characters are
1765 described in the locale definition file using “big-endian” notation for reasons of portability.
1766 There is no requirement that the internal representation in the computer memory be in this same
1767 order.

1768 One of the guidelines used for the development of this volume of IEEE Std 1003.1-200x is that
1769 characters outside the invariant part of the ISO/IEC 646:1991 standard should not be used in
1770 portable specifications. The backslash character is not in the invariant part; the number sign is,
1771 but with multiple representations: as a number sign, and as a pound sign. As far as general
1772 usage of these symbols, they are covered by the “grandfather clause”, but for newly defined
1773 interfaces, the WG15 POSIX working group has requested that POSIX provide alternate

1774 representations. Consequently, while the default escape character remains the backslash and the
1775 default comment character is the number sign, implementations are required to recognize
1776 alternative representations, identified in the applicable source file via the `<escape_char>` and
1777 `<comment_char>` keywords.

1778 A.7.3.1 *LC_CTYPE*

1779 The *LC_CTYPE* category is primarily used to define the encoding-independent aspects of a
1780 character set, such as character classification. In addition, certain encoding-dependent
1781 characteristics are also defined for an application via the *LC_CTYPE* category.
1782 IEEE Std 1003.1-200x does not mandate that the encoding used in the locale is the same as the
1783 one used by the application because an implementation may decide that it is advantageous to
1784 define locales in a system-wide encoding rather than having multiple, logically identical locales
1785 in different encodings, and to convert from the application encoding to the system-wide
1786 encoding on usage. Other implementations could require encoding-dependent locales.

1787 In either case, the *LC_CTYPE* attributes that are directly dependent on the encoding, such as
1788 `<mb_cur_max>` and the display width of characters, are not user-specifiable in a locale source
1789 and are consequently not defined as keywords.

1790 Implementations may define additional keywords or extend the *LC_CTYPE* mechanism to allow
1791 application-defined keywords.

1792 The text “The ellipsis specification shall only be valid within a single encoded character set” is
1793 present because it is possible to have a locale supported by multiple character encodings, as
1794 explained in the rationale for the Base Definitions volume of IEEE Std 1003.1-200x, Section 6.1,
1795 Portable Character Set. An example given there is of a possible Japanese-based locale supported
1796 by a mixture of the character sets JIS X 0201 Roman, JIS X 0208, and JIS X 0201 Katakana.
1797 Attempting to express a range of characters across these sets is not logical and the
1798 implementation is free to reject such attempts.

1799 As the *LC_CTYPE* character classes are based on the ISO C standard character class definition,
1800 the category does not support multi-character elements. For instance, the German character
1801 `<sharp-s>` is traditionally classified as a lowercase letter. There is no corresponding uppercase
1802 letter; in proper capitalization of German text, the `<sharp-s>` will be replaced by "SS"; that is, by
1803 two characters. This kind of conversion is outside the scope of the **toupper** and **tolower**
1804 keywords.

1805 Where IEEE Std 1003.1-200x specifies that only certain characters can be specified, as for the
1806 keywords **digit** and **xdigit**, the specified characters shall be from the portable character set, as
1807 shown. As an example, only the Arabic digits 0 through 9 are acceptable as digits.

1808 The character classes **digit**, **xdigit**, **lower**, **upper**, and **space** have a set of automatically included
1809 characters. These only need to be specified if the character values (that is, encoding) differs from
1810 the implementation default values. It is not possible to define a locale without these
1811 automatically included characters unless some implementation extension is used to prevent
1812 their inclusion. Such a definition would not be a proper superset of the C locale, and thus, it
1813 might not be possible for the standard utilities to be implemented as programs conforming to
1814 the ISO C standard.

1815 The definition of character class **digit** requires that only ten characters—the ones defining
1816 digits—can be specified; alternate digits (for example, Hindi or Kanji) cannot be specified here.
1817 However, the encoding may vary if an implementation supports more than one encoding.

1818 The definition of character class **xdigit** requires that the characters included in character class
1819 **digit** are included here also and allows for different symbols for the hexadecimal digits 10
1820 through 15.

1821 The inclusion of the **charclass** keyword satisfies the following requirement from the
1822 ISO POSIX-2: 1993 standard, Annex H.1:

1823 (3) *The LC_CTYPE (2.5.2.1) locale definition should be enhanced to allow user-specified additional*
1824 *character classes, similar in concept to the ISO C standard Multibyte Support Extension (MSE)*
1825 *is_wctype() function.*

1826 This keyword was previously included in The Open Group specifications and is now mandated
1827 in the Shell and Utilities volume of IEEE Std 1003.1-200x.

1828 The symbolic constant {CHARCLASS_NAME_MAX} was also adopted from The Open Group
1829 specifications. Application portability is enhanced by the use of symbolic constants.

1830 A.7.3.2 LC_COLLATE

1831 The rules governing collation depend to some extent on the use. At least five different levels of
1832 increasingly complex collation rules can be distinguished:

- 1833 1. *Byte/machine code order*: This is the historical collation order in the UNIX system and many
1834 proprietary operating systems. Collation is here performed character by character, without
1835 any regard to context. The primary virtue is that it usually is quite fast and also
1836 completely deterministic; it works well when the native machine collation sequence
1837 matches the user expectations.
- 1838 2. *Character order*: On this level, collation is also performed character by character, without
1839 regard to context. The order between characters is, however, not determined by the code
1840 values, but on the expectations by the user of the “correct” order between characters. In
1841 addition, such a (simple) collation order can specify that certain characters collate equally
1842 (for example, uppercase and lowercase letters).
- 1843 3. *String ordering*: On this level, entire strings are compared based on relatively
1844 straightforward rules. Several “passes” may be required to determine the order between
1845 two strings. Characters may be ignored in some passes, but not in others; the strings may
1846 be compared in different directions; and simple string substitutions may be performed
1847 before strings are compared. This level is best described as “dictionary” ordering; it is
1848 based on the spelling, not the pronunciation, or meaning, of the words.
- 1849 4. *Text search ordering*: This is a further refinement of the previous level, best described as
1850 “telephone book ordering”; some common homonyms (words spelled differently but with
1851 the same pronunciation) are collated together; numbers are collated as if they were spelled
1852 out, and so on.
- 1853 5. *Semantic-level ordering*: Words and strings are collated based on their meaning; entire words
1854 (such as “the”) are eliminated; the ordering is not deterministic. This usually requires
1855 special software and is highly dependent on the intended use.

1856 While the historical collation order formally is at level 1, for the English language it corresponds
1857 roughly to elements at level 2. The user expects to see the output from the *ls* utility sorted very
1858 much as it would be in a dictionary. While telephone book ordering would be an optimal goal
1859 for standard collation, this was ruled out as the order would be language-dependent.
1860 Furthermore, a requirement was that the order must be determined solely from the text string
1861 and the collation rules; no external information (for example, “pronunciation dictionaries”)
1862 could be required.

1863 As a result, the goal for the collation support is at level 3. This also matches the requirements for
1864 the Canadian collation order, as well as other, known collation requirements for alphabetic
1865 scripts. It specifically rules out collation based on pronunciation rules or based on semantic

1866 analysis of the text.

1867 The syntax for the *LC_COLLATE* category source meets the requirements for level 3 and has
1868 been verified to produce the correct result with examples based on French, Canadian, and
1869 Danish collation order. Because it supports multi-character collating elements, it is also capable
1870 of supporting collation in codesets where a character is expressed using non-spacing characters
1871 followed by the base character (such as the ISO/IEC 6937: 1994 standard).

1872 The directives that can be specified in an operand to the **order_start** keyword are based on the
1873 requirements specified in several proposed standards and in customary use. The following is a
1874 rephrasing of rules defined for “lexical ordering in English and French” by the Canadian
1875 Standards Association (the text in square brackets is rephrased):

- 1876 • Once special characters [punctuation] have been removed from original strings, the ordering
1877 is determined by scanning forwards (left to right) [disregarding case and diacriticals].
- 1878 • In case of equivalence, special characters are once again removed from original strings and
1879 the ordering is determined by scanning backwards (starting from the rightmost character of
1880 the string and back), character by character [disregarding case but considering diacriticals].
- 1881 • In case of repeated equivalence, special characters are removed again from original strings
1882 and the ordering is determined by scanning forwards, character by character [considering
1883 both case and diacriticals].
- 1884 • If there is still an ordering equivalence after the first three rules have been applied, then only
1885 special characters and the position they occupy in the string are considered to determine
1886 ordering. The string that has a special character in the lowest position comes first. If two
1887 strings have a special character in the same position, the character [with the lowest collation
1888 value] comes first. In case of equality, the other special characters are considered until there
1889 is a difference or until all special characters have been exhausted.

1890 It is estimated that this part of IEEE Std 1003.1-200x covers the requirements for all European
1891 languages, and no particular problems are anticipated with Slavic or Middle East character sets.

1892 The Far East (particularly Japanese/Chinese) collations are often based on contextual
1893 information and pronunciation rules (the same ideogram can have different meanings and
1894 different pronunciations). Such collation, in general, falls outside the desired goal of
1895 IEEE Std 1003.1-200x. There are, however, several other collation rules (stroke/radical or “most
1896 common pronunciation”) that can be supported with the mechanism described here.

1897 The character order is defined by the order in which characters and elements are specified
1898 between the **order_start** and **order_end** keywords. Weights assigned to the characters and
1899 elements define the collation sequence; in the absence of weights, the character order is also the
1900 collation sequence.

1901 The **position** keyword provides the capability to consider, in a compare, the relative position of
1902 characters not subject to **IGNORE**. As an example, consider the two strings "o-ring" and
1903 "or-ing". Assuming the hyphen is subject to **IGNORE** on the first pass, the two strings
1904 compare equal, and the position of the hyphen is immaterial. On second pass, all characters
1905 except the hyphen are subject to **IGNORE**, and in the normal case the two strings would again
1906 compare equal. By taking position into account, the first collates before the second.

1907 A.7.3.3 LC_MONETARY

1908 The currency symbol does not appear in LC_MONETARY because it is not defined in the C locale
 1909 of the ISO C standard.

1910 The ISO C standard limits the size of decimal points and thousands delimiters to single-byte
 1911 values. In locales based on multi-byte coded character sets, this cannot be enforced;
 1912 IEEE Std 1003.1-200x does not prohibit such characters, but makes the behavior unspecified (in
 1913 the text “In contexts where other standards ...”).

1914 The grouping specification is based on, but not identical to, the ISO C standard. The -1 signals
 1915 that no further grouping shall be performed; the equivalent of {CHAR_MAX} in the ISO C
 1916 standard.

1917 The text “the value is not available in the locale” is taken from the ISO C standard and is used
 1918 instead of the “unspecified” text in early proposals. There is no implication that omitting these
 1919 keywords or assigning them values of " " or -1 produces unspecified results; such omissions or
 1920 assignments eliminate the effects described for the keyword or produce zero-length strings, as
 1921 appropriate.

1922 The locale definition is an extension of the ISO C standard localeconv() specification. In
 1923 particular, rules on how currency_symbol is treated are extended to also cover int_curr_symbol,
 1924 and p_sep_by_space and n_sep_by_space have been augmented with the value 2, which places
 1925 a <space> between the sign and the symbol (if they are adjacent; otherwise, it should be treated
 1926 as a 0). The following table 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\$+

1939 The following is an example of the interpretation of the mon_grouping keyword. Assuming that
 1940 the value to be formatted is 123456789 and the mon_thousands_sep is ' ', then the following
 1941 table shows the result. The third column shows the equivalent string in the ISO C standard that
 1942 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"

1949 In these examples, the octal value of {CHAR_MAX} is 177.

1950 A.7.3.4 *LC_NUMERIC*

1951 See the rationale for *LC_MONETARY* for a description of the behavior of grouping.

1952 A.7.3.5 *LC_TIME*

1953 Although certain of the conversion specifications in the POSIX locale (such as the name of the
1954 month) are shown with initial capital letters, this need not be the case in other locales. Programs
1955 using these conversion specifications may need to adjust the capitalization if the output is going
1956 to be used at the beginning of a sentence.

1957 The *LC_TIME* descriptions of **abday**, **day**, **mon**, and **abmon** imply a Gregorian style calendar (7-
1958 day weeks, 12-month years, leap years, and so on). Formatting time strings for other types of
1959 calendars is outside the scope of IEEE Std 1003.1-200x.

1960 While the ISO 8601:2000 standard numbers the weekdays starting with Monday, historical
1961 practice is to use the Sunday as the first day. Rather than change the order and introduce
1962 potential confusion, the days must be specified beginning with Sunday; previous references to
1963 “first day” have been removed. Note also that the Shell and Utilities volume of
1964 IEEE Std 1003.1-200x *date* utility supports numbering compliant with the ISO 8601:2000
1965 standard.

1966 As specified under *date* in the Shell and Utilities volume of IEEE Std 1003.1-200x and *strftime()* in
1967 the System Interfaces volume of IEEE Std 1003.1-200x, the conversion specifications
1968 corresponding to the optional keywords consist of a modifier followed by a traditional
1969 conversion specification (for instance, %Ex). If the optional keywords are not supported by the
1970 implementation or are unspecified for the current locale, these modified conversion
1971 specifications are treated as the traditional conversion specifications. For example, assume the
1972 following keywords:

```
1973     alt_digits    "0th";"1st";"2nd";"3rd";"4th";"5th";\  
1974                "6th";"7th";"8th";"9th";"10th"  
1975     d_fmt       "The %Od day of %B in %Y"
```

1976 On July 4th 1776, the %x conversion specifications would result in "The 4th day of July
1977 in 1776", while on July 14th 1789 it would result in "The 14 day of July in 1789". It
1978 can be noted that the above example is for illustrative purposes only; the %O modifier is
1979 primarily intended to provide for Kanji or Hindi digits in *date* formats.

1980 The following is an example for Japan that supports the current plus last three Emperors and
1981 reverts to Western style numbering for years prior to the Meiji era. The example also allows for
1982 the custom of using a special name for the first year of an era instead of using 1. (The examples
1983 substitute romaji where kanji should be used.)

```
1984     era_d_fmt    "%EY%mgatsu%dnichi (%a)"  
1985     era         "+:2:1990/01/01:+*:Heisei:%EC%Eynen";\  
1986                "+:1:1989/01/08:1989/12/31:Heisei:%ECgannen";\  
1987                "+:2:1927/01/01:1989/01/07:Shouwa:%EC%Eynen";\  
1988                "+:1:1926/12/25:1926/12/31:Shouwa:%ECgannen";\  
1989                "+:2:1913/01/01:1926/12/24:Taishou:%EC%Eynen";\  
1990                "+:1:1912/07/30:1912/12/31:Taishou:%ECgannen";\  
1991                "+:2:1869/01/01:1912/07/29:Meiji:%EC%Eynen";\  
1992                "+:1:1868/09/08:1868/12/31:Meiji:%ECgannen";\  
1993                "-:1868:1868/09/07:-*:%Ey"
```

1994 Assuming that the current date is September 21, 1991, a request to *date* or *strftime()* would yield
 1995 the following results:

```
1996      %Ec - Heisei3nen9gatsu21nichi (Sat) 14:39:26
1997      %EC - Heisei
1998      %Ex - Heisei3nen9gatsu21nichi (Sat)
1999      %Ey - 3
2000      %EY - Heisei3nen
```

2001 Example era definitions for the Republic of China:

```
2002      era      "+:2:1913/01/01:+*:ChungHwaMingGuo:%EC%EyNen";\
2003      "+:1:1912/1/1:1912/12/31:ChungHwaMingGuo:%ECYuenNen";\
2004      "+:1:1911/12/31:-*:MingChien:%EC%EyNen"
```

2005 Example definitions for the Christian Era:

```
2006      era      "+:0:0001/01/01:+*:AD:%EC %Ey";\
2007      "+:1:-0001/12/31:-*:BC:%Ey %EC"
```

2008 A.7.3.6 LC_MESSAGES

2009 The **yesstr** and **nostr** locale keywords and the YESSTR and NOSTR *langinfo* items were formerly
 2010 used to match user affirmative and negative responses. In IEEE Std 1003.1-200x, the **yesexpr**,
 2011 **noexpr**, YESEXPR, and NOEXPR extended regular expressions have replaced them. |
 2012 Applications should use the general locale-based messaging facilities to issue prompting |
 2013 messages which include sample desired responses. |

2014 A.7.4 Locale Definition Grammar

2015 There is no additional rationale provided for this section.

2016 A.7.4.1 Locale Lexical Conventions

2017 There is no additional rationale provided for this section.

2018 A.7.4.2 Locale Grammar

2019 There is no additional rationale provided for this section.

2020 A.7.5 Locale Definition Example

2021 The following is an example of a locale definition file that could be used as input to the *localedef*
 2022 utility. It assumes that the utility is executed with the *-f* option, naming a *charmap* file with (at
 2023 least) the following content:

```

2024     CHARMAP
2025     <space>      \x20
2026     <dollar>     \x24
2027     <A>         \101
2028     <a>         \141
2029     <A-acute>   \346
2030     <a-acute>   \365
2031     <A-grave>  \300
2032     <a-grave>  \366
2033     <b>         \142
2034     <C>         \103
2035     <c>         \143
2036     <c-cedilla> \347
2037     <d>         \x64
2038     <H>         \110
2039     <h>         \150
2040     <eszet>    \xb7
2041     <s>         \x73
2042     <z>         \x7a
2043     END CHARMAP

```

2044 It should not be taken as complete or to represent any actual locale, but only to illustrate the
 2045 syntax.

```

2046     #
2047     LC_CTYPE
2048     lower  <a>;<b>;<c>;<c-cedilla>;<d>;...;<z>
2049     upper  A;B;C;Ç;...;Z
2050     space  \x20;\x09;\x0a;\x0b;\x0c;\x0d
2051     blank  \040;\011
2052     toupper (<a>,<A>);(b,B);(c,C);(ç,Ç);(d,D);(z,Z)
2053     END LC_CTYPE
2054     #
2055     LC_COLLATE
2056     #
2057     # The following example of collation is based on
2058     # Canadian standard Z243.4.1-1998, "Canadian Alphanumeric
2059     # Ordering Standard For Character sets of CSA Z234.4 Standard".
2060     # (Other parts of this example locale definition file do not
2061     # purport to relate to Canada, or to any other real culture.)
2062     # The proposed standard defines a 4-weight collation, such that
2063     # in the first pass, characters are compared without regard to
2064     # case or accents; in second pass, backwards compare without
2065     # regard to case; in the third pass, forward compare without
2066     # regard to diacriticals. In the 3 first passes, non-alphabetic
2067     # characters are ignored; in the fourth pass, only special
2068     # characters are considered, such that "The string that has a
2069     # special character in the lowest position comes first. If two
2070     # strings have a special character in the same position, the
2071     # collation value of the special character determines ordering.
2072     #
2073     # Only a subset of the character set is used here; mostly to
2074     # illustrate the set-up.

```

```

2075      #
2076      collating-symbol <NULL>
2077      collating-symbol <LOW_VALUE>
2078      collating-symbol <LOWER-CASE>
2079      collating-symbol <SUBSCRIPT-LOWER>
2080      collating-symbol <SUPERSCRIPT-LOWER>
2081      collating-symbol <UPPER-CASE>
2082      collating-symbol <NO-ACCENT>
2083      collating-symbol <PECULIAR>
2084      collating-symbol <LIGATURE>
2085      collating-symbol <ACUTE>
2086      collating-symbol <GRAVE>
2087      # Further collating-symbols follow.
2088      #
2089      # Properly, the standard does not include any multi-character
2090      # collating elements; the one below is added for completeness.
2091      #
2092      collating_element <ch> from "<c><h>"
2093      collating_element <CH> from "<C><H>"
2094      collating_element <Ch> from "<C><h>"
2095      #
2096      order_start forward;backward;forward;forward,position
2097      #
2098      # Collating symbols are specified first in the sequence to allocate
2099      # basic collation values to them, lower than that of any character.
2100      <NULL>
2101      <LOW_VALUE>
2102      <LOWER-CASE>
2103      <SUBSCRIPT-LOWER>
2104      <SUPERSCRIPT-LOWER>
2105      <UPPER-CASE>
2106      <NO-ACCENT>
2107      <PECULIAR>
2108      <LIGATURE>
2109      <ACUTE>
2110      <GRAVE>
2111      <RING-ABOVE>
2112      <DIAERESIS>
2113      <TILDE>
2114      # Further collating symbols are given a basic collating value here.
2115      #
2116      # Here follow special characters.
2117      <space>          IGNORE;IGNORE;IGNORE;<space>
2118      # Other special characters follow here.
2119      #
2120      # Here follow the regular characters.
2121      <a>              <a>;<NO-ACCENT>;<LOWER-CASE>;IGNORE
2122      <A>              <a>;<NO-ACCENT>;<UPPER-CASE>;IGNORE
2123      <a-acute>        <a>;<ACUTE>;<LOWER-CASE>;IGNORE
2124      <A-acute>        <a>;<ACUTE>;<UPPER-CASE>;IGNORE
2125      <a-grave>        <a>;<GRAVE>;<LOWER-CASE>;IGNORE
2126      <A-grave>        <a>;<GRAVE>;<UPPER-CASE>;IGNORE

```

```

2127 <ae>      "<a><e>" ; "<LIGATURE><LIGATURE>" ; \
2128          "<LOWER-CASE><LOWER-CASE>" ; IGNORE
2129 <AE>      "<a><e>" ; "<LIGATURE><LIGATURE>" ; \
2130          "<UPPER-CASE><UPPER-CASE>" ; IGNORE
2131 <b>        <b> ; <NO-ACCENT> ; <LOWER-CASE> ; IGNORE
2132 <B>        <b> ; <NO-ACCENT> ; <UPPER-CASE> ; IGNORE
2133 <c>        <c> ; <NO-ACCENT> ; <LOWER-CASE> ; IGNORE
2134 <C>        <c> ; <NO-ACCENT> ; <UPPER-CASE> ; IGNORE
2135 <ch>       <ch> ; <NO-ACCENT> ; <LOWER-CASE> ; IGNORE
2136 <Ch>       <ch> ; <NO-ACCENT> ; <PECULIAR> ; IGNORE
2137 <CH>       <ch> ; <NO-ACCENT> ; <UPPER-CASE> ; IGNORE
2138 #
2139 # As an example, the strings "Bach" and "bach" could be encoded (for
2140 # compare purposes) as:
2141 # "Bach"   <b> ; <a> ; <ch> ; <LOW_VALUE> ; <NO-ACCENT> ; <NO-ACCENT> ; \
2142 #          <NO-ACCENT> ; <LOW_VALUE> ; <UPPER-CASE> ; <LOWER-CASE> ; \
2143 #          <LOWER-CASE> ; <NULL>
2144 # "bach"   <b> ; <a> ; <ch> ; <LOW_VALUE> ; <NO-ACCENT> ; <NO-ACCENT> ; \
2145 #          <NO-ACCENT> ; <LOW_VALUE> ; <LOWER-CASE> ; <LOWER-CASE> ; \
2146 #          <LOWER-CASE> ; <NULL>
2147 #
2148 # The two strings are equal in pass 1 and 2, but differ in pass 3.
2149 #
2150 # Further characters follow.
2151 #
2152 UNDEFINED  IGNORE ; IGNORE ; IGNORE ; IGNORE
2153 #
2154 order_end
2155 #
2156 END LC_COLLATE
2157 #
2158 LC_MONETARY
2159 int_curr_symbol      "USD "
2160 currency_symbol     "$"
2161 mon_decimal_point    "."
2162 mon_grouping         3;0
2163 positive_sign        ""
2164 negative_sign        "- "
2165 p_cs_precedes        1
2166 n_sign_posn          0
2167 END LC_MONETARY
2168 #
2169 LC_NUMERIC
2170 copy "US_en.ASCII"
2171 END LC_NUMERIC
2172 #
2173 LC_TIME
2174 abday   "Sun" ; "Mon" ; "Tue" ; "Wed" ; "Thu" ; "Fri" ; "Sat"
2175 #
2176 day     "Sunday" ; "Monday" ; "Tuesday" ; "Wednesday" ; \
2177         "Thursday" ; "Friday" ; "Saturday"
2178 #

```

```

2179     abmon    "Jan"; "Feb"; "Mar"; "Apr"; "May"; "Jun"; \
2180             "Jul"; "Aug"; "Sep"; "Oct"; "Nov"; "Dec"
2181     #
2182     mon      "January"; "February"; "March"; "April"; \
2183             "May"; "June"; "July"; "August"; "September"; \
2184             "October"; "November"; "December"
2185     #
2186     d_t_fmt  "%a %b %d %T %Z %Y\n"
2187     END LC_TIME
2188     #
2189     LC_MESSAGES
2190     yesexpr  "^[yY][[:alpha:]]*" | (OK) "
2191     #
2192     noexpr   "^[nN][[:alpha:]]*"
2193     END LC_MESSAGES

```

2194 **A.8 Environment Variables**

2195 **A.8.1 Environment Variable Definition**

2196 The variable *environ* is not intended to be declared in any header, but rather to be declared by the
 2197 user for accessing the array of strings that is the environment. This is the traditional usage of the
 2198 symbol. Putting it into a header could break some programs that use the symbol for their own
 2199 purposes.

2200 The decision to restrict conforming systems to the use of digits, uppercase letters, and
 2201 underscores for environment variable names allows applications to use lowercase letters in their
 2202 environment variable names without conflicting with any conforming system.

2203 In addition to the obvious conflict with the shell syntax for positional parameter substitution,
 2204 some historical applications (including some shells) exclude names with leading digits from the
 2205 environment.

2206 **A.8.2 Internationalization Variables**

2207 The text about locale implies that any utilities written in standard C and conforming to
 2208 IEEE Std 1003.1-200x must issue the following call:

```
2209     setlocale(LC_ALL, "")
```

2210 If this were omitted, the ISO C standard specifies that the C locale would be used.

2211 If any of the environment variables are invalid, it makes sense to default to an implementation-
 2212 defined, consistent locale environment. It is more confusing for a user to have partial settings
 2213 occur in case of a mistake. All utilities would then behave in one language/cultural
 2214 environment. Furthermore, it provides a way of forcing the whole environment to be the
 2215 implementation-defined default. Disastrous results could occur if a pipeline of utilities partially
 2216 uses the environment variables in different ways. In this case, it would be appropriate for
 2217 utilities that use *LANG* and related variables to exit with an error if any of the variables are
 2218 invalid. For example, users typing individual commands at a terminal might want *date* to work if
 2219 *LC_MONETARY* is invalid as long as *LC_TIME* is valid. Since these are conflicting reasonable
 2220 alternatives, IEEE Std 1003.1-200x leaves the results unspecified if the locale environment
 2221 variables would not produce a complete locale matching the specification of the user.

2222 The locale settings of individual categories cannot be truly independent and still guarantee
 2223 correct results. For example, when collating two strings, characters must first be extracted from
 2224 each string (governed by *LC_CTYPE*) before being mapped to collating elements (governed by
 2225 *LC_COLLATE*) for comparison. That is, if *LC_CTYPE* is causing parsing according to the rules of
 2226 a large, multi-byte code set (potentially returning 20 000 or more distinct character codeset
 2227 values), but *LC_COLLATE* is set to handle only an 8-bit codeset with 256 distinct characters,
 2228 meaningful results are obviously impossible.

2229 The *LC_MESSAGES* variable affects the language of messages generated by the standard
 2230 utilities.

2231 The description of the environment variable names starting with the characters “LC_”
 2232 acknowledges the fact that the interfaces presented may be extended as new international
 2233 functionality is required. In the ISO C standard, names preceded by “LC_” are reserved in the
 2234 name space for future categories.

2235 To avoid name clashes, new categories and environment variables are divided into two
 2236 classifications: *implementation-independent* and *implementation-defined*.

2237 Implementation-independent names will have the following format:

2238 `LC_NAME`

2239 where *NAME* is the name of the new category and environment variable. Capital letters must be
 2240 used for implementation-independent names.

2241 Implementation-defined names must be in lowercase letters, as below:

2242 `LC_name`

2243 A.8.3 Other Environment Variables

2244 The quoted form of the timezone variable allows timezone names of the form UTC+1 (or any
 2245 name that contains the character plus ('+'), the character minus ('-'), or digits), which may be
 2246 appropriate for countries that do not have an official timezone name. It would be coded as
 2247 <UTC+1>+1<UTC+2>, which would cause *std* to have a value of UTC+1 and *dst* a value of
 2248 UTC+2, each with a length of 5 characters. This does not appear to conflict with any existing
 2249 usage. The characters '<' and '>' were chosen for quoting because they are easier to parse
 2250 visually than a quoting character that does not provide some sense of bracketing (and in a string
 2251 like this, such bracketing is helpful). They were also chosen because they do not need special
 2252 treatment when assigning to the *TZ* variable. Users are often confused by embedding quotes in a
 2253 string. Because '<' and '>' are meaningful to the shell, the whole string would have to be
 2254 quoted, but that is easily explained. (Parentheses would have presented the same problems.)
 2255 Although the '>' symbol could have been permitted in the string by either escaping it or
 2256 doubling it, it seemed of little value to require that. This could be provided as an extension if
 2257 there was a need. Timezone names of this new form lead to a requirement that the value of
 2258 `{_POSIX_TZNAME_MAX}` change from 3 to 6.

2259 COLUMNS, LINES

2260 The default value for the number of column positions, *COLUMNS*, and screen height, *LINES*, are
 2261 unspecified because historical implementations use different methods to determine values
 2262 corresponding to the size of the screen in which the utility is run. This size is typically known to
 2263 the implementation through the value of *TERM*, or by more elaborate methods such as
 2264 extensions to the *stty* utility or knowledge of how the user is dynamically resizing windows on a
 2265 bit-mapped display terminal. Users should not need to set these variables in the environment
 2266 unless there is a specific reason to override the default behavior of the implementation, such as

2267 to display data in an area arbitrarily smaller than the terminal or window. Values for these
 2268 variables that are not decimal integers greater than zero are implicitly undefined values; it is
 2269 unnecessary to enumerate all of the possible values outside of the acceptable set.

2270 **PATH**

2271 Many historical implementations of the Bourne shell do not interpret a trailing colon to represent
 2272 the current working directory and are thus non-conforming. The C Shell and the KornShell
 2273 conform to IEEE Std 1003.1-200x on this point. The usual name of dot may also be used to refer
 2274 to the current working directory.

2275 Many implementations historically have used a default value of **/bin** and **/usr/bin** for the *PATH*
 2276 variable. IEEE Std 1003.1-200x does not mandate this default path be identical to that retrieved
 2277 from *getconf _CS_PATH* because it is likely that the standardized utilities may be provided in
 2278 another directory separate from the directories used by some historical applications.

2279 **LOGNAME**

2280 In most implementations, the value of such a variable is easily forged, so security-critical
 2281 applications should rely on other means of determining user identity. *LOGNAME* is required to
 2282 be constructed from the portable filename character set for reasons of interchange. No diagnostic
 2283 condition is specified for violating this rule, and no requirement for enforcement exists. The
 2284 intent of the requirement is that if extended characters are used, the “guarantee” of portability
 2285 implied by a standard is void.

2286 **SHELL**

2287 The *SHELL* variable names the preferred shell of the user; it is a guide to applications. There is
 2288 no direct requirement that that shell conform to IEEE Std 1003.1-200x; that decision should rest
 2289 with the user. It is the intention of the standard developers that alternative shells be permitted, if
 2290 the user chooses to develop or acquire one. An operating system that builds its shell into the
 2291 “kernel” in such a manner that alternative shells would be impossible does not conform to the
 2292 spirit of IEEE Std 1003.1-200x.

2293 **CHANGE HISTORY**

2294 **Issue 6**

2295 Changed format of *TZ* field to allow for the quoted form as defined in previous
 2296 versions of the ISO POSIX-1 standard.

2297 **A.9 Regular Expressions**

2298 Rather than repeating the description of REs for each utility supporting REs, the standard
 2299 developers preferred a common, comprehensive description of regular expressions in one place.
 2300 The most common behavior is described here, and exceptions or extensions to this are
 2301 documented for the respective utilities, as appropriate.

2302 The BRE corresponds to the *ed* or historical *grep* type, and the ERE corresponds to the historical
 2303 *egrep* type (now *grep -E*).

2304 The text is based on the *ed* description and substantially modified, primarily to aid developers
 2305 and others in the understanding of the capabilities and limitations of REs. Much of this was
 2306 influenced by internationalization requirements.

2307 It should be noted that the definitions in this section do not cover the *tr* utility; the *tr* syntax does
2308 not employ REs.

2309 The specification of REs is particularly important to internationalization because pattern
2310 matching operations are very basic operations in business and other operations. The syntax and
2311 rules of REs are intended to be as intuitive as possible to make them easy to understand and use.
2312 The historical rules and behavior do not provide that capability to non-English language users,
2313 and do not provide the necessary support for commonly used characters and language
2314 constructs. It was necessary to provide extensions to the historical RE syntax and rules to
2315 accommodate other languages.

2316 As they are limited to bracket expressions, the rationale for these modifications is in the Base
2317 Definitions volume of IEEE Std 1003.1-200x, Section 9.3.5, RE Bracket Expression.

2318 A.9.1 Regular Expression Definitions

2319 It is possible to determine what strings correspond to subexpressions by recursively applying
2320 the leftmost longest rule to each subexpression, but only with the proviso that the overall match
2321 is leftmost longest. For example, matching "`\(ac*\)c*d[ac]*\1`" against *acdacaaa* matches
2322 *acdacaaa* (with `\1=a`); simply matching the longest match for "`\(ac*\)`" would yield `\1=ac`, but
2323 the overall match would be smaller (*acdac*). Conceptually, the implementation must examine
2324 every possible match and among those that yield the leftmost longest total matches, pick the one
2325 that does the longest match for the leftmost subexpression, and so on. Note that this means that
2326 matching by subexpressions is context-dependent: a subexpression within a larger RE may
2327 match a different string from the one it would match as an independent RE, and two instances of
2328 the same subexpression within the same larger RE may match different lengths even in similar
2329 sequences of characters. For example, in the ERE "`(a.*b)(a.*b)`", the two identical
2330 subexpressions would match four and six characters, respectively, of *accbacccb*.

2331 The definition of *single character* has been expanded to include also collating elements consisting
2332 of two or more characters; this expansion is applicable only when a bracket expression is
2333 included in the BRE or ERE. An example of such a collating element may be the Dutch *ij*, which
2334 collates as a 'Y'. In some encodings, a ligature "i with j" exists as a character and would
2335 represent a single-character collating element. In another encoding, no such ligature exists, and
2336 the two-character sequence *ij* is defined as a multi-character collating element. Outside brackets,
2337 the *ij* is treated as a two-character RE and matches the same characters in a string. Historically, a
2338 bracket expression only matched a single character. The ISO POSIX-2:1993 standard required
2339 bracket expressions like "`^[[:lower:]]`" to match multi-character collating elements such as
2340 "*ij*". However, this requirement led to behavior that many users did not expect and that could
2341 not feasibly be mimicked in user code, and it was rarely if ever implemented correctly. The
2342 current standard leaves it unspecified whether a bracket expression matches a multi-character
2343 collating element, allowing both historical and ISO POSIX-2:1993 standard implementations to
2344 conform.

2345 Also, in the current standard, it is unspecified whether character class expressions like
2346 "`[:lower:]`" can include multi-character collating elements like "*ij*"; hence
2347 "`[[:lower:]]`" can match "*ij*", and "`^[[:lower:]]`" can fail to match "*ij*". Common
2348 practice is for a character class expression to match a collating element if it matches the collating
2349 element's first character.

2350 **A.9.2 Regular Expression General Requirements**

2351 The definition of which sequence is matched when several are possible is based on the leftmost-
2352 longest rule historically used by deterministic recognizers. This rule is easier to define and
2353 describe, and arguably more useful, than the first-match rule historically used by non-
2354 deterministic recognizers. It is thought that dependencies on the choice of rule are rare; carefully
2355 contrived examples are needed to demonstrate the difference.

2356 A formal expression of the leftmost-longest rule is:

2357 The search is performed as if all possible suffixes of the string were tested for a prefix
2358 matching the pattern; the longest suffix containing a matching prefix is chosen, and the
2359 longest possible matching prefix of the chosen suffix is identified as the matching sequence.

2360 Historically, most RE implementations only match lines, not strings. However, that is more an
2361 effect of the usage than of an inherent feature of REs themselves. Consequently,
2362 IEEE Std 1003.1-200x does not regard <newline>s as special; they are ordinary characters, and
2363 both a period and a non-matching list can match them. Those utilities (like *grep*) that do not
2364 allow <newline>s to match are responsible for eliminating any <newline> from strings before
2365 matching against the RE. The *regcomp()* function, however, can provide support for such
2366 processing without violating the rules of this section.

2367 Some implementations of *egrep* have had very limited flexibility in handling complex EREs.
2368 IEEE Std 1003.1-200x does not attempt to define the complexity of a BRE or ERE, but does place
2369 a lower limit on it—any RE must be handled, as long as it can be expressed in 256 bytes or less.
2370 (Of course, this does not place an upper limit on the implementation.) There are historical
2371 programs using a non-deterministic-recognizer implementation that should have no difficulty
2372 with this limit. It is possible that a good approach would be to attempt to use the faster, but
2373 more limited, deterministic recognizer for simple expressions and to fall back on the non-
2374 deterministic recognizer for those expressions requiring it. Non-deterministic implementations
2375 must be careful to observe the rules on which match is chosen; the longest match, not the first
2376 match, starting at a given character is used.

2377 The term *invalid* highlights a difference between this section and some others:
2378 IEEE Std 1003.1-200x frequently avoids mandating of errors for syntax violations because they
2379 can be used by implementors to trigger extensions. However, the authors of the
2380 internationalization features of REs wanted to mandate errors for certain conditions to identify
2381 usage problems or non-portable constructs. These are identified within this rationale as
2382 appropriate. The remaining syntax violations have been left implicitly or explicitly undefined.
2383 For example, the BRE construct " $\{1,2,3\}$ " does not comply with the grammar. A
2384 conforming application cannot rely on it producing an error nor matching the literal characters
2385 " $\{1,2,3\}$ ". The term "undefined" was used in favor of "unspecified" because many of the
2386 situations are considered errors on some implementations, and the standard developers
2387 considered that consistency throughout the section was preferable to mixing undefined and
2388 unspecified.

2389 **A.9.3 Basic Regular Expressions**

2390 There is no additional rationale provided for this section.

2391 *A.9.3.1 BREs Matching a Single Character or Collating Element*

2392 There is no additional rationale provided for this section.

2393 *A.9.3.2 BRE Ordinary Characters*

2394 There is no additional rationale provided for this section.

2395 *A.9.3.3 BRE Special Characters*

2396 There is no additional rationale provided for this section.

2397 *A.9.3.4 Periods in BREs*

2398 There is no additional rationale provided for this section.

2399 *A.9.3.5 RE Bracket Expression*

2400 Range expressions are, historically, an integral part of REs. However, the requirements of
 2401 “natural language behavior” and portability do conflict. In the POSIX locale, ranges must be
 2402 treated according to the collating sequence and include such characters that fall within the range
 2403 based on that collating sequence, regardless of character values. In other locales, ranges have
 2404 unspecified behavior.

2405 Some historical implementations allow range expressions where the ending range point of one
 2406 range is also the starting point of the next (for instance, "[a-m-o]"). This behavior should not
 2407 be permitted, but to avoid breaking historical implementations, it is now *undefined* whether it is a
 2408 valid expression and how it should be interpreted.

2409 Current practice in *awk* and *lex* is to accept escape sequences in bracket expressions as per the
 2410 Base Definitions volume of IEEE Std 1003.1-200x, Table 5-1, Escape Sequences and Associated
 2411 Actions, while the normal ERE behavior is to regard such a sequence as consisting of two
 2412 characters. Allowing the *awk/lex* behavior in EREs would change the normal behavior in an
 2413 unacceptable way; it is expected that *awk* and *lex* will decode escape sequences in EREs before
 2414 passing them to *regcomp()* or comparable routines. Each utility describes the escape sequences it
 2415 accepts as an exception to the rules in this section; the list is not the same, for historical reasons.

2416 As noted previously, the new syntax and rules have been added to accommodate other
 2417 languages than English. The remainder of this section describes the rationale for these
 2418 modifications.

2419 In the POSIX locale, a regular expression that starts with a range expression matches a set of
 2420 strings that are contiguously sorted, but this is not necessarily true in other locales. For example,
 2421 a French locale might have the following behavior:

```
2422 $ ls
2423 alpha Alpha estimé ESTIMÉ été eurêka
2424 $ ls [a-e]*
2425 alpha Alpha estimé eurêka
```

2426 Such disagreements between matching and contiguous sorting are unavoidable because POSIX
 2427 sorting cannot be implemented in terms of a deterministic finite-state automaton (DFA), but
 2428 range expressions by design are implementable in terms of DFAs.

2429 Historical implementations used native character order to interpret range expressions. The
 2430 ISO POSIX-2:1993 standard instead required collating element order (CEO): the order that
 2431 collating elements were specified between the **order_start** and **order_end** keywords in the
 2432 *LC_COLLATE* category of the current locale. CEO had some advantages in portability over the
 2433 native character order, but it also had some disadvantages:

- 2434 • CEO could not feasibly be mimicked in user code, leading to inconsistencies between POSIX
 2435 matchers and matchers in popular user programs like Emacs, *ksh*, and Perl.
- 2436 • CEO caused range expressions to match accented and capitalized letters contrary to many
 2437 users' expectations. For example, "[a-e]" typically matched both 'E' and 'á' but neither
 2438 'A' nor 'é'.
- 2439 • CEO was not consistent across implementations. In practice, CEO was often less portable
 2440 than native character order. For example, it was common for the CEOs of two
 2441 implementation-supplied locales to disagree, even if both locales were named "da_DK".

2442 Because of these problems, some implementations of regular expressions continued to use
 2443 native character order. Others used the collation sequence, which is more consistent with sorting
 2444 than either CEO or native order, but which departs further from the traditional POSIX semantics
 2445 because it generally requires "[a-e]" to match either 'A' or 'E' but not both. As a result of
 2446 this kind of implementation variation, programmers who wanted to write portable regular
 2447 expressions could not rely on the ISO POSIX-2:1993 standard guarantees in practice.

2448 While revising the standard, lengthy consideration was given to proposals to attack this problem
 2449 by adding an API for querying the CEO to allow user-mode matchers, but none of these
 2450 proposals had implementation experience and none achieved consensus. Leaving the standard
 2451 alone was also considered, but rejected due to the problems described above.

2452 The current standard leaves unspecified the behavior of a range expression outside the POSIX
 2453 locale. This makes it clearer that conforming applications should avoid range expressions
 2454 outside the POSIX locale, and it allows implementations and compatible user-mode matchers to
 2455 interpret range expressions using native order, CEO, collation sequence, or other, more
 2456 advanced techniques.

2457 The ISO POSIX-2:1993 standard required "[b-a]" to be an invalid expression in the POSIX
 2458 locale, but this requirement has been relaxed in this version of the standard so that "[b-a]" can
 2459 instead be treated as a valid expression that does not match any string.

2460 A.9.3.6 BREs Matching Multiple Characters

2461 The limit of nine back-references to subexpressions in the RE is based on the use of a single-digit
 2462 identifier; increasing this to multiple digits would break historical applications. This does not
 2463 imply that only nine subexpressions are allowed in REs. The following is a valid BRE with ten
 2464 subexpressions:

```
2465 \(\(\(ab\) *c\) *d\)\(ef\) *\ (gh)\{2\}\(ij\) *\ (kl)\ *(mn)\ *(op)\ *(qr)\ *
```

2466 The standard developers regarded the common historical behavior, which supported "\n*", but
 2467 not "\n\{min,max\}", "\(\.\.\)*", or "\(\.\.\)\{min,max\}", as a non-intentional
 2468 result of a specific implementation, and they supported both duplication and interval
 2469 expressions following subexpressions and back-references.

2470 The changes to the processing of the back-reference expression remove an unspecified or
 2471 ambiguous behavior in the Shell and Utilities volume of IEEE Std 1003.1-200x, aligning it with
 2472 the requirements specified for the *regcomp()* expression, and is the result of PASC Interpretation
 2473 1003.2-92 #43 submitted for the ISO POSIX-2:1993 standard.

2474 A.9.3.7 BRE Precedence

2475 There is no additional rationale provided for this section.

2476 A.9.3.8 BRE Expression Anchoring

2477 Often, the dollar sign is viewed as matching the ending <newline> in text files. This is not
2478 strictly true; the <newline> is typically eliminated from the strings to be matched, and the dollar
2479 sign matches the terminating null character.

2480 The ability of '^', '\$', and '*' to be non-special in certain circumstances may be confusing to
2481 some programmers, but this situation was changed only in a minor way from historical practice
2482 to avoid breaking many historical scripts. Some consideration was given to making the use of
2483 the anchoring characters undefined if not escaped and not at the beginning or end of strings.
2484 This would cause a number of historical BREs, such as "2^10", "\$HOME", and "\$1.35", that
2485 relied on the characters being treated literally, to become invalid.

2486 However, one relatively uncommon case was changed to allow an extension used on some
2487 implementations. Historically, the BREs "^foo" and "\(^foo\)" did not match the same
2488 string, despite the general rule that subexpressions and entire BREs match the same strings. To
2489 increase consensus, IEEE Std 1003.1-200x has allowed an extension on some implementations to
2490 treat these two cases in the same way by declaring that anchoring *may* occur at the beginning or
2491 end of a subexpression. Therefore, portable BREs that require a literal circumflex at the
2492 beginning or a dollar sign at the end of a subexpression must escape them. Note that a BRE such
2493 as "a\(^bc\)" will either match "a^bc" or nothing on different systems under the rules.

2494 ERE anchoring has been different from BRE anchoring in all historical systems. An unescaped
2495 anchor character has never matched its literal counterpart outside a bracket expression. Some
2496 implementations treated "foo\$bar" as a valid expression that never matched anything; others
2497 treated it as invalid. IEEE Std 1003.1-200x mandates the former, valid unmatched behavior.

2498 Some implementations have extended the BRE syntax to add alternation. For example, the
2499 subexpression "\ (foo\$|bar\)" would match either "foo" at the end of the string or "bar"
2500 anywhere. The extension is triggered by the use of the undefined "\|" sequence. Because the
2501 BRE is undefined for portable scripts, the extending system is free to make other assumptions,
2502 such that the '\$' represents the end-of-line anchor in the middle of a subexpression. If it were
2503 not for the extension, the '\$' would match a literal dollar sign under the rules.

2504 A.9.4 Extended Regular Expressions

2505 As with BREs, the standard developers decided to make the interpretation of escaped ordinary
2506 characters undefined.

2507 The right parenthesis is not listed as an ERE special character because it is only special in the
2508 context of a preceding left parenthesis. If found without a preceding left parenthesis, the right
2509 parenthesis has no special meaning.

2510 The *interval expression*, "{m,n}", has been added to EREs. Historically, the interval expression
2511 has only been supported in some ERE implementations. The standard developers estimated that
2512 the addition of interval expressions to EREs would not decrease consensus and would also make
2513 BREs more of a subset of EREs than in many historical implementations.

2514 It was suggested that, in addition to interval expressions, back-references ('\n') should also be
2515 added to EREs. This was rejected by the standard developers as likely to decrease consensus.

2516 In historical implementations, multiple duplication symbols are usually interpreted from left to
2517 right and treated as additive. As an example, "a+b" matches zero or more instances of 'a'
2518 followed by a 'b'. In IEEE Std 1003.1-200x, multiple duplication symbols are undefined; that is,

- 2519 they cannot be relied upon for conforming applications. One reason for this is to provide some
2520 scope for future enhancements.
- 2521 The precedence of operations differs between EREs and those in *lex*; in *lex*, for historical reasons,
2522 interval expressions have a lower precedence than concatenation.
- 2523 **A.9.4.1 EREs Matching a Single Character or Collating Element**
- 2524 There is no additional rationale provided for this section.
- 2525 **A.9.4.2 ERE Ordinary Characters**
- 2526 There is no additional rationale provided for this section.
- 2527 **A.9.4.3 ERE Special Characters**
- 2528 There is no additional rationale provided for this section.
- 2529 **A.9.4.4 Periods in EREs**
- 2530 There is no additional rationale provided for this section.
- 2531 **A.9.4.5 ERE Bracket Expression**
- 2532 There is no additional rationale provided for this section.
- 2533 **A.9.4.6 EREs Matching Multiple Characters**
- 2534 There is no additional rationale provided for this section.
- 2535 **A.9.4.7 ERE Alternation**
- 2536 There is no additional rationale provided for this section.
- 2537 **A.9.4.8 ERE Precedence**
- 2538 There is no additional rationale provided for this section.
- 2539 **A.9.4.9 ERE Expression Anchoring**
- 2540 There is no additional rationale provided for this section.
- 2541 **A.9.5 Regular Expression Grammar**
- 2542 The grammars are intended to represent the range of acceptable syntaxes available to
2543 conforming applications. There are instances in the text where undefined constructs are
2544 described; as explained previously, these allow implementation extensions. There is no intended
2545 requirement that an implementation extension must somehow fit into the grammars shown
2546 here.
- 2547 The BRE grammar does not permit L_ANCHOR or R_ANCHOR inside "`\(`" and "`\)`" (which
2548 implies that '`^`' and '`$`' are ordinary characters). This reflects the semantic limits on the
2549 application, as noted in the Base Definitions volume of IEEE Std 1003.1-200x, Section 9.3.8, BRE
2550 Expression Anchoring. Implementations are permitted to extend the language to interpret '`^`'
2551 and '`$`' as anchors in these locations, and as such, conforming applications cannot use
2552 unescaped '`^`' and '`$`' in positions inside "`\(`" and "`\)`" that might be interpreted as anchors.
- 2553 The ERE grammar does not permit several constructs that the Base Definitions volume of
2554 IEEE Std 1003.1-200x, Section 9.4.2, ERE Ordinary Characters and the Base Definitions volume of

2555 IEEE Std 1003.1-200x, Section 9.4.3, ERE Special Characters specify as having undefined results:
 2556 • ORD_CHAR preceded by '\'
 2557 • *ERE_dupl_symbol*(s) appearing first in an ERE, or immediately following '|', '^', or '('
 2558 • '{' not part of a valid *ERE_dupl_symbol*
 2559 • '|' appearing first or last in an ERE, or immediately following '|' or '(', or immediately
 2560 preceding ')'
 2561 Implementations are permitted to extend the language to allow these. Conforming applications |
 2562 cannot use such constructs. |

2563 A.9.5.1 BRE/ERE Grammar Lexical Conventions

2564 There is no additional rationale provided for this section.

2565 A.9.5.2 RE and Bracket Expression Grammar

2566 The removal of the *Back_open_paren Back_close_paren* option from the *nondupl_RE* specification is
 2567 the result of PASC Interpretation 1003.2-92 #43 submitted for the ISO POSIX-2: 1993 standard.
 2568 Although the grammar required support for null subexpressions, this section does not describe
 2569 the meaning of, and historical practice did not support, this construct.

2570 A.9.5.3 ERE Grammar

2571 There is no additional rationale provided for this section.

2572 A.10 Directory Structure and Devices

2573 A.10.1 Directory Structure and Files

2574 A description of the historical **/usr/tmp** was omitted, removing any concept of differences in
 2575 emphasis between the / and **/usr** directories. The descriptions of **/bin**, **/usr/bin**, **/lib**, and **/usr/lib**
 2576 were omitted because they are not useful for applications. In an early draft, a distinction was
 2577 made between system and application directory usage, but this was not found to be useful.

2578 The directories / and **/dev** are included because the notion of a hierarchical directory structure is
 2579 key to other information presented elsewhere in IEEE Std 1003.1-200x. In early drafts, it was
 2580 argued that special devices and temporary files could conceivably be handled without a
 2581 directory structure on some implementations. For example, the system could treat the characters
 2582 `"/tmp"` as a special token that would store files using some non-POSIX file system structure.
 2583 This notion was rejected by the standard developers, who required that all the files in this
 2584 section be implemented via POSIX file systems.

2585 The **/tmp** directory is retained in IEEE Std 1003.1-200x to accommodate historical applications
 2586 that assume its availability. Implementations are encouraged to provide suitable directory
 2587 names in the environment variable *TMPDIR* and applications are encouraged to use the contents
 2588 of *TMPDIR* for creating temporary files.

2589 The standard files **/dev/null** and **/dev/tty** are required to be both readable and writable to allow
 2590 applications to have the intended historical access to these files.

2591 The standard file **/dev/console** has been added for alignment with the Single UNIX Specification.

2592 **A.10.2 Output Devices and Terminal Types**

2593 There is no additional rationale provided for this section.

2594 **A.11 General Terminal Interface**

2595 If the implementation does not support this interface on any device types, it should behave as if
 2596 it were being used on a device that is not a terminal device (in most cases *errno* will be set to
 2597 [ENOTTY] on return from functions defined by this interface). This is based on the fact that
 2598 many applications are written to run both interactively and in some non-interactive mode, and
 2599 they adapt themselves at runtime. Requiring that they all be modified to test an environment
 2600 variable to determine whether they should try to adapt is unnecessary. On a system that
 2601 provides no general terminal interface, providing all the entry points as stubs that return
 2602 [ENOTTY] (or an equivalent, as appropriate) has the same effect and requires no changes to the
 2603 application.

2604 Although the needs of both interface implementors and application developers were addressed
 2605 throughout IEEE Std 1003.1-200x, this section pays more attention to the needs of the latter. This
 2606 is because, while many aspects of the programming interface can be hidden from the user by the
 2607 application developer, the terminal interface is usually a large part of the user interface.
 2608 Although to some extent the application developer can build missing features or work around
 2609 inappropriate ones, the difficulties of doing that are greater in the terminal interface than
 2610 elsewhere. For example, efficiency prohibits the average program from interpreting every
 2611 character passing through it in order to simulate character erase, line kill, and so on. These
 2612 functions should usually be done by the operating system, possibly at the interrupt level.

2613 The *tc**() functions were introduced as a way of avoiding the problems inherent in the
 2614 traditional *ioctl*() function and in variants of it that were proposed. For example, *tcsetattr*()
 2615 is specified in place of the use of the TCSETA *ioctl*() command function. This allows specification
 2616 of all the arguments in a manner consistent with the ISO C standard unlike the varying third
 2617 argument of *ioctl*(), which is sometimes a pointer (to any of many different types) and
 2618 sometimes an **int**.

2619 The advantages of this new method include:

- 2620 • It allows strict type checking.
- 2621 • The direction of transfer of control data is explicit.
- 2622 • Portable capabilities are clearly identified.
- 2623 • The need for a general interface routine is avoided.
- 2624 • Size of the argument is well-defined (there is only one type).

2625 The disadvantages include:

- 2626 • No historical implementation used the new method.
- 2627 • There are many small routines instead of one general-purpose one.
- 2628 • The historical parallel with *fcntl*() is broken.

2629 The issue of modem control was excluded from IEEE Std 1003.1-200x on the grounds that:

- 2630 • It was concerned with setting and control of hardware timers.
- 2631 • The appropriate timers and settings vary widely internationally.

2632 • Feedback from European computer manufacturers indicated that this facility was not
2633 consistent with European needs and that specification of such a facility was not a
2634 requirement for portability.

2635 **A.11.1 Interface Characteristics**

2636 *A.11.1.1 Opening a Terminal Device File*

2637 There is no additional rationale provided for this section.

2638 *A.11.1.2 Process Groups*

2639 There is a potential race when the members of the foreground process group on a terminal leave
2640 that process group, either by exit or by changing process groups. After the last process exits the
2641 process group, but before the foreground process group ID of the terminal is changed (usually
2642 by a job-control shell), it would be possible for a new process to be created with its process ID
2643 equal to the terminal's foreground process group ID. That process might then become the
2644 process group leader and accidentally be placed into the foreground on a terminal that was not
2645 necessarily its controlling terminal. As a result of this problem, the controlling terminal is
2646 defined to not have a foreground process group during this time.

2647 The cases where a controlling terminal has no foreground process group occur when all
2648 processes in the foreground process group either terminate and are waited for or join other
2649 process groups via *setpgid()* or *setsid()*. If the process group leader terminates, this is the first
2650 case described; if it leaves the process group via *setpgid()*, this is the second case described (a
2651 process group leader cannot successfully call *setsid()*). When one of those cases causes a
2652 controlling terminal to have no foreground process group, it has two visible effects on
2653 applications. The first is the value returned by *tcgetpgrp()*. The second (which occurs only in the
2654 case where the process group leader terminates) is the sending of signals in response to special
2655 input characters. The intent of IEEE Std 1003.1-200x is that no process group be wrongly
2656 identified as the foreground process group by *tcgetpgrp()* or unintentionally receive signals
2657 because of placement into the foreground.

2658 In 4.3 BSD, the old process group ID continues to be used to identify the foreground process
2659 group and is returned by the function equivalent to *tcgetpgrp()*. In that implementation it is
2660 possible for a newly created process to be assigned the same value as a process ID and then form
2661 a new process group with the same value as a process group ID. The result is that the new
2662 process group would receive signals from this terminal for no apparent reason, and
2663 IEEE Std 1003.1-200x precludes this by forbidding a process group from entering the foreground
2664 in this way. It would be more direct to place part of the requirement made by the last sentence
2665 under *fork()*, but there is no convenient way for that section to refer to the value that *tcgetpgrp()*
2666 returns, since in this case there is no process group and thus no process group ID.

2667 One possibility for a conforming implementation is to behave similarly to 4.3 BSD, but to
2668 prevent this reuse of the ID, probably in the implementation of *fork()*, as long as it is in use by
2669 the terminal.

2670 Another possibility is to recognize when the last process stops using the terminal's foreground
2671 process group ID, which is when the process group lifetime ends, and to change the terminal's
2672 foreground process group ID to a reserved value that is never used as a process ID or process
2673 group ID. (See the definition of *process group lifetime* in the definitions section.) The process ID
2674 can then be reserved until the terminal has another foreground process group.

2675 The 4.3 BSD implementation permits the leader (and only member) of the foreground process
2676 group to leave the process group by calling the equivalent of *setpgid()* and to later return,
2677 expecting to return to the foreground. There are no known application needs for this behavior,

2678 and IEEE Std 1003.1-200x neither requires nor forbids it (except that it is forbidden for session
2679 leaders) by leaving it unspecified.

2680 A.11.1.3 *The Controlling Terminal*

2681 IEEE Std 1003.1-200x does not specify a mechanism by which to allocate a controlling terminal.
2682 This is normally done by a system utility (such as *getty*) and is considered an administrative
2683 feature outside the scope of IEEE Std 1003.1-200x.

2684 Historical implementations allocate controlling terminals on certain *open()* calls. Since *open()* is
2685 part of POSIX.1, its behavior had to be dealt with. The traditional behavior is not required
2686 because it is not very straightforward or flexible for either implementations or applications.
2687 However, because of its prevalence, it was not practical to disallow this behavior either. Thus, a
2688 mechanism was standardized to ensure portable, predictable behavior in *open()*.

2689 Some historical implementations deallocate a controlling terminal on the last system-wide close.
2690 This behavior is neither required nor prohibited. Even on implementations that do provide this
2691 behavior, applications generally cannot depend on it due to its system-wide nature.

2692 A.11.1.4 *Terminal Access Control*

2693 The access controls described in this section apply only to a process that is accessing its
2694 controlling terminal. A process accessing a terminal that is not its controlling terminal is
2695 effectively treated the same as a member of the foreground process group. While this may seem
2696 unintuitive, note that these controls are for the purpose of job control, not security, and job
2697 control relates only to a process' controlling terminal. Normal file access permissions handle
2698 security.

2699 If the process calling *read()* or *write()* is in a background process group that is orphaned, it is not
2700 desirable to stop the process group, as it is no longer under the control of a job control shell that
2701 could put it into foreground again. Accordingly, calls to *read()* or *write()* functions by such
2702 processes receive an immediate error return. This is different from 4.2 BSD, which kills orphaned
2703 processes that receive terminal stop signals.

2704 The foreground/background/orphaned process group check performed by the terminal driver
2705 must be repeatedly performed until the calling process moves into the foreground or until the
2706 process group of the calling process becomes orphaned. That is, when the terminal driver
2707 determines that the calling process is in the background and should receive a job control signal,
2708 it sends the appropriate signal (SIGTTIN or SIGTTOU) to every process in the process group of
2709 the calling process and then it allows the calling process to immediately receive the signal. The
2710 latter is typically performed by blocking the process so that the signal is immediately noticed.
2711 Note, however, that after the process finishes receiving the signal and control is returned to the
2712 driver, the terminal driver must reexecute the foreground/background/orphaned process group
2713 check. The process may still be in the background, either because it was continued in the
2714 background by a job-control shell, or because it caught the signal and did nothing.

2715 The terminal driver repeatedly performs the foreground/background/orphaned process group
2716 checks whenever a process is about to access the terminal. In the case of *write()* or the control
2717 *tc*()* functions, the check is performed at the entry of the function. In the case of *read()*, the check
2718 is performed not only at the entry of the function, but also after blocking the process to wait for
2719 input characters (if necessary). That is, once the driver has determined that the process calling
2720 the *read()* function is in the foreground, it attempts to retrieve characters from the input queue. If
2721 the queue is empty, it blocks the process waiting for characters. When characters are available
2722 and control is returned to the driver, the terminal driver must return to the repeated
2723 foreground/background/orphaned process group check again. The process may have moved
2724 from the foreground to the background while it was blocked waiting for input characters.

2725 A.11.1.5 *Input Processing and Reading Data*

2726 There is no additional rationale provided for this section.

2727 A.11.1.6 *Canonical Mode Input Processing*

2728 The term *character* is intended here. ERASE should erase the last character, not the last byte. In
2729 the case of multi-byte characters, these two may be different.

2730 4.3 BSD has a WERASE character that erases the last “word” typed (but not any preceding
2731 <blank>s or <tab>s). A word is defined as a sequence of non-<blank>s, with <tab>s counted as
2732 <blank>s. Like ERASE, WERASE does not erase beyond the beginning of the line. This
2733 WERASE feature has not been specified in POSIX.1 because it is difficult to define in the
2734 international environment. It is only useful for languages where words are delimited by
2735 <blank>s. In some ideographic languages, such as Japanese and Chinese, words are not
2736 delimited at all. The WERASE character should presumably take one back to the beginning of a
2737 sentence in those cases; practically, this means it would not get much use for those languages.

2738 It should be noted that there is a possible inherent deadlock if the application and
2739 implementation conflict on the value of MAX_CANON. With ICANON set (if IXOFF is
2740 enabled) and more than MAX_CANON characters transmitted without a <linefeed>,
2741 transmission will be stopped, the <linefeed> (or <carriage-return> when ICRLF is set) will never
2742 arrive, and the *read()* will never be satisfied.

2743 An application should not set IXOFF if it is using canonical mode unless it knows that (even in
2744 the face of a transmission error) the conditions described previously cannot be met or unless it is
2745 prepared to deal with the possible deadlock in some other way, such as timeouts.

2746 It should also be noted that this can be made to happen in non-canonical mode if the trigger
2747 value for sending IXOFF is less than VMIN and VTIME is zero.

2748 A.11.1.7 *Non-Canonical Mode Input Processing*

2749 Some points to note about MIN and TIME:

- 2750 1. The interactions of MIN and TIME are not symmetric. For example, when MIN>0 and
2751 TIME=0, TIME has no effect. However, in the opposite case where MIN=0 and TIME>0,
2752 both MIN and TIME play a role in that MIN is satisfied with the receipt of a single
2753 character.
- 2754 2. Also note that in case A (MIN>0, TIME>0), TIME represents an inter-character timer, while
2755 in case C (MIN=0, TIME>0), TIME represents a read timer.

2756 These two points highlight the dual purpose of the MIN/TIME feature. Cases A and B, where
2757 MIN>0, exist to handle burst-mode activity (for example, file transfer programs) where a
2758 program would like to process at least MIN characters at a time. In case A, the inter-character
2759 timer is activated by a user as a safety measure; in case B, it is turned off.

2760 Cases C and D exist to handle single-character timed transfers. These cases are readily adaptable
2761 to screen-based applications that need to know if a character is present in the input queue before
2762 refreshing the screen. In case C, the read is timed; in case D, it is not.

2763 Another important note is that MIN is always just a minimum. It does not denote a record
2764 length. That is, if a program does a read of 20 bytes, MIN is 10, and 25 characters are present, 20
2765 characters shall be returned to the user. In the special case of MIN=0, this still applies: if more
2766 than one character is available, they all will be returned immediately.

2767 **A.11.1.8 Writing Data and Output Processing**

2768 There is no additional rationale provided for this section.

2769 **A.11.1.9 Special Characters**

2770 There is no additional rationale provided for this section.

2771 **A.11.1.10 Modem Disconnect**

2772 There is no additional rationale provided for this section.

2773 **A.11.1.11 Closing a Terminal Device File**

2774 IEEE Std 1003.1-200x does not specify that a *close()* on a terminal device file include the
2775 equivalent of a call to *tcfow(fd,TCOON)*.

2776 An implementation that discards output at the time *close()* is called after reporting the return
2777 value to the *write()* call that data was written does not conform with IEEE Std 1003.1-200x. An
2778 application has functions such as *tcdrain()*, *tcfush()*, and *tcfow()* available to obtain the detailed
2779 behavior it requires with respect to flushing of output.

2780 At the time of the last close on a terminal device, an application relinquishes any ability to exert
2781 flow control via *tcfow()*.

2782 **A.11.2 Parameters that Can be Set**2783 **A.11.2.1 The termios Structure**

2784 This structure is part of an interface that, in general, retains the historic grouping of flags.
2785 Although a more optimal structure for implementations may be possible, the degree of change
2786 to applications would be significantly larger.

2787 **A.11.2.2 Input Modes**

2788 Some historical implementations treated a long break as multiple events, as many as one per
2789 character time. The wording in POSIX.1 explicitly prohibits this.

2790 Although the ISTRIP flag is normally superfluous with today's terminal hardware and software,
2791 it is historically supported. Therefore, applications may be using ISTRIP, and there is no
2792 technical problem with supporting this flag. Also, applications may wish to receive only 7-bit
2793 input bytes and may not be connected directly to the hardware terminal device (for example,
2794 when a connection traverses a network).

2795 Also, there is no requirement in general that the terminal device ensures that high-order bits
2796 beyond the specified character size are cleared. ISTRIP provides this function for 7-bit
2797 characters, which are common.

2798 In dealing with multi-byte characters, the consequences of a parity error in such a character, or in
2799 an escape sequence affecting the current character set, are beyond the scope of POSIX.1 and are
2800 best dealt with by the application processing the multi-byte characters.

2801 *A.11.2.3 Output Modes*

2802 POSIX.1 does not describe postprocessing of output to a terminal or detailed control of that from |
2803 a conforming application. (That is, translation of <newline> to <carriage-return> followed by |
2804 <linefeed> or <tab> processing.) There is nothing that a conforming application should do to its |
2805 output for a terminal because that would require knowledge of the operation of the terminal. It
2806 is the responsibility of the operating system to provide postprocessing appropriate to the output
2807 device, whether it is a terminal or some other type of device.

2808 Extensions to POSIX.1 to control the type of postprocessing already exist and are expected to
2809 continue into the future. The control of these features is primarily to adjust the interface between
2810 the system and the terminal device so the output appears on the display correctly. This should
2811 be set up before use by any application.

2812 In general, both the input and output modes should not be set absolutely, but rather modified
2813 from the inherited state.

2814 *A.11.2.4 Control Modes*

2815 This section could be misread that the symbol “CSIZE” is a title in the **termios** *c_cflag* field .
2816 Although it does serve that function, it is also a required symbol, as a literal reading of POSIX.1
2817 (and the caveats about typography) would indicate.

2818 *A.11.2.5 Local Modes*

2819 Non-canonical mode is provided to allow fast bursts of input to be read efficiently while still
2820 allowing single-character input.

2821 The ECHONL function historically has been in many implementations. Since there seems to be
2822 no technical problem with supporting ECHONL, it is included in POSIX.1 to increase consensus.

2823 The alternate behavior possible when ECHOK or ECHOE are specified with ICANON is
2824 permitted as a compromise depending on what the actual terminal hardware can do. Erasing
2825 characters and lines is preferred, but is not always possible.

2826 *A.11.2.6 Special Control Characters*

2827 Permitting VMIN and VTIME to overlap with VEOF and VEOL was a compromise for historical
2828 implementations. Only when backwards-compatibility of object code is a serious concern to an
2829 implementor should an implementation continue this practice. Correct applications that work
2830 with the overlap (at the source level) should also work if it is not present, but not the reverse.

2831 **A.12 Utility Conventions**2832 **A.12.1 Utility Argument Syntax**

2833 The standard developers considered that recent trends toward diluting the SYNOPSIS sections
2834 of historical reference pages to the equivalent of:

2835 `command [options][operands]`

2836 were a disservice to the reader. Therefore, considerable effort was placed into rigorous
2837 definitions of all the command line arguments and their interrelationships. The relationships
2838 depicted in the synopses are normative parts of IEEE Std 1003.1-200x; this information is
2839 sometimes repeated in textual form, but that is only for clarity within context.

2840 The use of “undefined” for conflicting argument usage and for repeated usage of the same |
2841 option is meant to prevent conforming applications from using conflicting arguments or |
2842 repeated options unless specifically allowed (as is the case with *ls*, which allows simultaneous,
2843 repeated use of the *-C*, *-l*, and *-1* options). Many historical implementations will tolerate this
2844 usage, choosing either the first or the last applicable argument. This tolerance can continue, but |
2845 conforming applications cannot rely upon it. (Other implementations may choose to print usage |
2846 messages instead.)

2847 The use of “undefined” for conflicting argument usage also allows an implementation to make
2848 reasonable extensions to utilities where the implementor considers mutually-exclusive options
2849 according to IEEE Std 1003.1-200x to have a sensible meaning and result.

2850 IEEE Std 1003.1-200x does not define the result of a command when an option-argument or
2851 operand is not followed by ellipses and the application specifies more than one of that option-
2852 argument or operand. This allows an implementation to define valid (although non-standard)
2853 behavior for the utility when more than one such option or operand is specified.

2854 Allowing <blank>s after an option (that is, placing an option and its option-argument into
2855 separate argument strings) when IEEE Std 1003.1-200x does not require it encourages portability
2856 of users, while still preserving backwards-compatibility of scripts. Inserting <blank>s between
2857 the option and the option-argument is preferred; however, historical usage has not been
2858 consistent in this area; therefore, <blank>s are required to be handled by all implementations,
2859 but implementations are also allowed to handle the historical syntax. Another justification for
2860 selecting the multiple-argument method was that the single-argument case is inherently
2861 ambiguous when the option-argument can legitimately be a null string.

2862 IEEE Std 1003.1-200x explicitly states that digits are permitted as operands and option-
2863 arguments. The lower and upper bounds for the values of the numbers used for operands and
2864 option-arguments were derived from the ISO C standard values for {LONG_MIN} and
2865 {LONG_MAX}. The requirement on the standard utilities is that numbers in the specified range
2866 do not cause a syntax error, although the specification of a number need not be semantically
2867 correct for a particular operand or option-argument of a utility. For example, the specification of:

2868 `dd obs=3000000000`

2869 would yield undefined behavior for the application and could be a syntax error because the
2870 number 3 000 000 000 is outside of the range $-2\ 147\ 483\ 647$ to $+2\ 147\ 483\ 647$. On the other hand:

2871 `dd obs=2000000000`

2872 may cause some error, such as “blocksize too large”, rather than a syntax error.

2873 **A.12.2 Utility Syntax Guidelines**

2874 This section is based on the rules listed in the SVID. It was included for two reasons:

- 2875 1. The individual utility descriptions in the Shell and Utilities volume of
2876 IEEE Std 1003.1-200x, Chapter 4, Utilities needed a set of common (although not universal)
2877 actions on which they could anchor their descriptions of option and operand syntax. Most
2878 of the standard utilities actually do use these guidelines, and many of their historical
2879 implementations use the *getopt()* function for their parsing. Therefore, it was simpler to
2880 cite the rules and merely identify exceptions.
- 2881 2. Writers of conforming applications need suggested guidelines if the POSIX community is
2882 to avoid the chaos of historical UNIX system command syntax.

2883 It is recommended that all *future* utilities and applications use these guidelines to enhance “user
2884 portability”. The fact that some historical utilities could not be changed (to avoid breaking
2885 historical applications) should not deter this future goal.

2886 The voluntary nature of the guidelines is highlighted by repeated uses of the word *should*
2887 throughout. This usage should not be misinterpreted to imply that utilities that claim
2888 conformance in their OPTIONS sections do not always conform.

2889 Guidelines 1 and 2 are offered as guidance for locales using Latin alphabets. No
2890 recommendations are made by IEEE Std 1003.1-200x concerning utility naming in other locales.

2891 In the Shell and Utilities volume of IEEE Std 1003.1-200x, Section 2.9.1, Simple Commands, it is
2892 further stated that a command used in the Shell Command Language cannot be named with a
2893 trailing colon.

2894 Guideline 3 was changed to allow alphanumeric characters (letters and digits) from the character
2895 set to allow compatibility with historical usage. Historical practice allows the use of digits
2896 wherever practical, and there are no portability issues that would prohibit the use of digits. In
2897 fact, from an internationalization viewpoint, digits (being non-language-dependent) are
2898 preferable over letters (a *-2* is intuitively self-explanatory to any user, while in the *-f filename* the
2899 letter ‘f’ is a mnemonic aid only to speakers of Latin-based languages where “filename”
2900 happens to translate to a word that begins with ‘f’). Since guideline 3 still retains the word
2901 “single”, multi-digit options are not allowed. Instances of historical utilities that used them have
2902 been marked obsolescent, with the numbers being changed from option names to option-
2903 arguments.

2904 It was difficult to achieve a satisfactory solution to the problem of name space in option
2905 characters. When the standard developers desired to extend the historical *cc* utility to accept
2906 ISO C standard programs, they found that all of the portable alphabet was already in use by
2907 various vendors. Thus, they had to devise a new name, *c89*, rather than something like *cc -X*.
2908 There were suggestions that implementors be restricted to providing extensions through various
2909 means (such as using a plus sign as the option delimiter or using option characters outside the
2910 alphanumeric set) that would reserve all of the remaining alphanumeric characters for future
2911 POSIX standards. These approaches were resisted because they lacked the historical style of
2912 UNIX systems. Furthermore, if a vendor-provided option should become commonly used in the
2913 industry, it would be a candidate for standardization. It would be desirable to standardize such a
2914 feature using historical practice for the syntax (the semantics can be standardized with any
2915 syntax). This would not be possible if the syntax was one reserved for the vendor. However,
2916 since the standardization process may lead to minor changes in the semantics, it may prove to be
2917 better for a vendor to use a syntax that will not be affected by standardization.

2918 Guideline 8 includes the concept of comma-separated lists in a single argument. It is up to the
2919 utility to parse such a list itself because *getopt()* just returns the single string. This situation was

2920 retained so that certain historical utilities would not violate the guidelines. Applications
2921 preparing for international use should be aware of an occasional problem with comma-
2922 separated lists: in some locales, the comma is used as the radix character. Thus, if an application
2923 is preparing operands for a utility that expects a comma-separated lists, it should avoid
2924 generating non-integer values through one of the means that is influenced by setting the
2925 *LC_NUMERIC* variable (such as *awk*, *bc*, *printf*, or *printf()*).

2926 Applications calling any utility with a first operand starting with '-' should usually specify --,
2927 as indicated by Guideline 10, to mark the end of the options. This is true even if the SYNOPSIS in
2928 the Shell and Utilities volume of IEEE Std 1003.1-200x does not specify any options;
2929 implementations may provide options as extensions to the Shell and Utilities volume of
2930 IEEE Std 1003.1-200x. The standard utilities that do not support Guideline 10 indicate that fact in
2931 the OPTIONS section of the utility description.

2932 Guideline 11 was modified to clarify that the order of different options should not matter
2933 relative to one another. However, the order of repeated options that also have option-arguments
2934 may be significant; therefore, such options are required to be interpreted in the order that they
2935 are specified. The *make* utility is an instance of a historical utility that uses repeated options
2936 in which the order is significant. Multiple files are specified by giving multiple instances of the -f
2937 option; for example:

```
2938     make -f common_header -f specific_rules target
```

2939 Guideline 13 does not imply that all of the standard utilities automatically accept the operand
2940 '-' to mean standard input or output, nor does it specify the actions of the utility upon
2941 encountering multiple '-' operands. It simply says that, by default, '-' operands are not used
2942 for other purposes in the file reading or writing (but not when using *stat*, *unlink*, *touch*, and so on)
2943 utilities. All information concerning actual treatment of the '-' operand is found in the
2944 individual utility sections.

2945 An area of concern was that as implementations mature, implementation-defined utilities and
2946 implementation-defined utility options will result. The idea was expressed that there needed to
2947 be a standard way, say an environment variable or some such mechanism, to identify
2948 implementation-defined utilities separately from standard utilities that may have the same
2949 name. It was decided that there already exist several ways of dealing with this situation and that
2950 it is outside of the POSIX.2 scope to attempt to standardize in the area of non-standard items. A
2951 method that exists on some historical implementations is the use of the so-called */local/bin* or
2952 */usr/local/bin* directory to separate local or additional copies or versions of utilities. Another
2953 method that is also used is to isolate utilities into completely separate domains. Still another
2954 method to ensure that the desired utility is being used is to request the utility by its full
2955 pathname. There are many approaches to this situation; the examples given above serve to
2956 illustrate that there is more than one.

2957 **A.13 Headers**2958 **A.13.1 Format of Entries**

2959 Each header reference page has a common layout of sections describing the interface. This layout
2960 is similar to the manual page or “man” page format shipped with most UNIX systems, and each
2961 header has sections describing the SYNOPSIS and DESCRIPTION. These are the two sections
2962 that relate to conformance.

2963 Additional sections are informative, and add considerable information for the application
2964 developer. APPLICATION USAGE sections provide additional caveats, issues, and
2965 recommendations to the developer. RATIONALE sections give additional information on the
2966 decisions made in defining the interface.

2967 FUTURE DIRECTIONS sections act as pointers to related work that may impact the interface in
2968 the future, and often cautions the developer to architect the code to account for a change in this
2969 area. Note that a future directions statement should not be taken as a commitment to adopt a
2970 feature or interface in the future.

2971 The CHANGE HISTORY section describes when the interface was introduced, and how it has
2972 changed.

2973 Option labels and margin markings in the page can be useful in guiding the application |
2974 developer.

2975 / *Rationale (Informative)*

2976 **Part B:**

2977 **System Interfaces**

2978 *The Open Group*
2979 *The Institute of Electrical and Electronics Engineers, Inc.*

Rationale for System Interfaces

2980

2981 **B.1 Introduction**

2982 **B.1.1 Scope**

2983 Refer to Section A.1.1 (on page 3293).

2984 **B.1.2 Conformance**

2985 Refer to Section A.2 (on page 3299).

2986 **B.1.3 Normative References**

2987 There is no additional rationale provided for this section. |

2988 **B.1.4 Change History**

2989 The change history is provided as an informative section, to track changes from previous issues |
2990 of IEEE Std 1003.1-200x. |

2991 The following sections describe changes made to the System Interfaces volume of |
2992 IEEE Std 1003.1-200x since Issue 5 of the base document. The CHANGE HISTORY section for |
2993 each entry details the technical changes that have been made to that entry since Issue 5. |
2994 Changes between earlier issues of the base document and Issue 5 are not included. |

2995 The change history between Issue 5 and Issue 6 also lists the changes since the |
2996 ISO POSIX-1: 1996 standard. |

2997 **Changes from Issue 5 to Issue 6 (IEEE Std 1003.1-200x)**

2998 The following list summarizes the major changes that were made in the System Interfaces |
2999 volume of IEEE Std 1003.1-200x from Issue 5 to Issue 6: |

- 3000 • This volume of IEEE Std 1003.1-200x is extensively revised so it can be both an IEEE POSIX |
3001 Standard and an Open Group Technical Standard. |
- 3002 • The POSIX System Interfaces requirements incorporate support of FIPS 151-2. |
- 3003 • The POSIX System Interfaces requirements are updated to align with some features of the |
3004 Single UNIX Specification. |
- 3005 • A RATIONALE section is added to each reference page. |
- 3006 • Networking interfaces from the XNS, Issue 5.2 specification are incorporated. |
- 3007 • IEEE Std 1003.1d-1999 is incorporated. |
- 3008 • IEEE Std 1003.1j-2000 is incorporated. |
- 3009 • IEEE Std 1003.1q-2000 is incorporated. |
- 3010 • IEEE P1003.1a draft standard is incorporated. |

- 3011 • Existing functionality is aligned with the ISO/IEC 9899: 1999 standard.
- 3012 • New functionality from the ISO/IEC 9899: 1999 standard is incorporated.
- 3013 • IEEE PASC Interpretations are applied.
- 3014 • The Open Group corrigenda and resolutions are applied.

3015 **New Features in Issue 6**

3016 The functions first introduced in Issue 6 (over the Issue 5 Base document) are listed in the table
 3017 below:

3018
 3019

3020
 3021
 3022
 3023
 3024
 3025
 3026
 3027
 3028
 3029
 3030
 3031
 3032
 3033
 3034
 3035
 3036
 3037
 3038
 3039
 3040
 3041
 3042
 3043
 3044
 3045
 3046
 3047
 3048
 3049
 3050
 3051
 3052
 3053
 3054
 3055
 3056

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

3057
3058
3059
3060
3061
3062
3063
3064
3065
3066
3067
3068
3069
3070
3071
3072
3073
3074
3075
3076
3077
3078
3079
3080
3081
3082
3083
3084
3085
3086
3087
3088
3089
3090
3091
3092
3093
3094
3095
3096
3097
3098
3099
3100
3101
3102

New Functions in Issue 6

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

3103
3104
3105
3106
3107
3108
3109
3110
3111
3112
3113
3114
3115
3116
3117
3118
3119
3120
3121
3122
3123
3124
3125
3126
3127
3128
3129
3130
3131
3132
3133
3134
3135
3136
3137
3138
3139
3140
3141
3142
3143

New Functions in Issue 6

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

3144 The following new headers are introduced in Issue 6:

3145

3146

3147

3148

3149

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

3150 The following table lists the functions and symbols from the XSI extension. These are new since
 3151 the ISO POSIX-1: 1996 standard.

3152
 3153

New XSI Functions and Symbols in Issue 6				
3154	<i>_longjmp()</i>	<i>getcontext()</i>	<i>msgget()</i>	<i>setutxent()</i>
3155	<i>_setjmp()</i>	<i>getdate()</i>	<i>msgrcv()</i>	<i>shmat()</i>
3156	<i>_tolower()</i>	<i>getgrent()</i>	<i>msgsnd()</i>	<i>shmctl()</i>
3157	<i>_toupper()</i>	<i>gethostid()</i>	<i>nftw()</i>	<i>shmdt()</i>
3158	<i>a64l()</i>	<i>getitimer()</i>	<i>nice()</i>	<i>shmget()</i>
3159	<i>basename()</i>	<i>getpgid()</i>	<i>nl_langinfo()</i>	<i>sigaltstack()</i>
3160	<i>bcmp()</i>	<i>getpmsg()</i>	<i>rand48()</i>	<i>sighold()</i>
3161	<i>bcopy()</i>	<i>getpriority()</i>	<i>openlog()</i>	<i>sigignore()</i>
3162	<i>bzero()</i>	<i>getpwent()</i>	<i>poll()</i>	<i>siginterrupt()</i>
3163	<i>catclose()</i>	<i>getrlimit()</i>	<i>pread()</i>	<i>sigpause()</i>
3164	<i>catgets()</i>	<i>getrusage()</i>	<i>pthread_attr_getguardsize()</i>	<i>sigrelse()</i>
3165	<i>catopen()</i>	<i>getsid()</i>	<i>pthread_attr_setguardsize()</i>	<i>sigset()</i>
3166	<i>closelog()</i>	<i>getsubopt()</i>	<i>pthread_attr_setstack()</i>	<i>srand48()</i>
3167	<i>crypt()</i>	<i>gettimeofday()</i>	<i>pthread_getconcurrency()</i>	<i>srandom()</i>
3168	<i>daylight</i>	<i>getutxent()</i>	<i>pthread_mutexattr_gettype()</i>	<i>statvfs()</i>
3169	<i>dbm_clearerr()</i>	<i>getutxid()</i>	<i>pthread_mutexattr_settype()</i>	<i>strcasecmp()</i>
3170	<i>dbm_close()</i>	<i>getutxline()</i>	<i>pthread_rwlockattr_init()</i>	<i>strdup()</i>
3171	<i>dbm_delete()</i>	<i>getwd()</i>	<i>pthread_rwlockattr_setpshared()</i>	<i>strfmon()</i>
3172	<i>dbm_error()</i>	<i>grantpt()</i>	<i>pthread_setconcurrency()</i>	<i>strncasecmp()</i>
3173	<i>dbm_fetch()</i>	<i>hcreate()</i>	<i>ptsname()</i>	<i>strptime()</i>
3174	<i>dbm_firstkey()</i>	<i>hdestroy()</i>	<i>putenv()</i>	<i>swab()</i>
3175	<i>dbm_nextkey()</i>	<i>hsearch()</i>	<i>pututxline()</i>	<i>swapcontext()</i>
3176	<i>dbm_open()</i>	<i>iconv()</i>	<i>pwrite()</i>	<i>sync()</i>
3177	<i>dbm_store()</i>	<i>iconv_close()</i>	<i>random()</i>	<i>syslog()</i>
3178	<i>dirname()</i>	<i>iconv_open()</i>	<i>readv()</i>	<i>tcgetsid()</i>
3179	<i>dlclose()</i>	<i>index()</i>	<i>realpath()</i>	<i>tdelete()</i>
3180	<i>dlderror()</i>	<i>initstate()</i>	<i>remque()</i>	<i>telldir()</i>
3181	<i>dlopen()</i>	<i>insque()</i>	<i>rindex()</i>	<i>tempnam()</i>
3182	<i>dlsym()</i>	<i>isascii()</i>	<i>seed48()</i>	<i>tfind()</i>
3183	<i>drand48()</i>	<i>jrand48()</i>	<i>seekdir()</i>	<i>timezone</i>
3184	<i>ecvt()</i>	<i>killpg()</i>	<i>semctl()</i>	<i>toascii()</i>
3185	<i>encrypt()</i>	<i>l64a()</i>	<i>semget()</i>	<i>truncate()</i>
3186	<i>endgrent()</i>	<i>lchown()</i>	<i>semop()</i>	<i>tsearch()</i>
3187	<i>endpwent()</i>	<i>lcong48()</i>	<i>setcontext()</i>	<i>twalk()</i>
3188	<i>endutxent()</i>	<i>lfind()</i>	<i>setgrent()</i>	<i>ulimit()</i>
3189	<i>erand48()</i>	<i>lockf()</i>	<i>setitimer()</i>	<i>unlockpt()</i>
3190	<i>fchdir()</i>	<i>lrand48()</i>	<i>setkey()</i>	<i>utimes()</i>
3191	<i>fcvt()</i>	<i>lsearch()</i>	<i>setlogmask()</i>	<i>waitid()</i>
3192	<i>ffs()</i>	<i>makecontext()</i>	<i>setpggrp()</i>	<i>wcswcs()</i>
3193	<i>fntmsg()</i>	<i>memccpy()</i>	<i>setpriority()</i>	<i>wcswidth()</i>
3194	<i>fstatvfs()</i>	<i>mknod()</i>	<i>setpwent()</i>	<i>wcwidth()</i>
3195	<i>ftime()</i>	<i>mkstemp()</i>	<i>setregid()</i>	<i>writev()</i>
3196	<i>ftok()</i>	<i>mktemp()</i>	<i>setreuid()</i>	
3197	<i>ftw()</i>	<i>mrnd48()</i>	<i>setrlimit()</i>	
3198	<i>gcvt()</i>	<i>msgctl()</i>	<i>setstate()</i>	

3199 The following table lists the headers from the XSI extension. These are new since the
 3200 ISO POSIX-1: 1996 standard.

3201
3202
3203
3204
3205
3206
3207
3208
3209
3210
3211
3212

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

3213 **B.1.5 Terminology**

3214 Refer to Section A.1.4 (on page 3295).

3215 **B.1.6 Definitions**

3216 Refer to Section A.3 (on page 3302).

3217 **B.1.7 Relationship to Other Formal Standards**

3218 There is no additional rationale provided for this section.

3219 **B.1.8 Portability**

3220 Refer to Section A.1.5 (on page 3298).

3221 *B.1.8.1 Codes*

3222 Refer to Section A.1.5.1 (on page 3298).

3223 **B.1.9 Format of Entries**

3224 Each system interface reference page has a common layout of sections describing the interface.
3225 This layout is similar to the manual page or “man” page format shipped with most UNIX
3226 systems, and each header has sections describing the SYNOPSIS, DESCRIPTION, RETURN
3227 VALUE, and ERRORS. These are the four sections that relate to conformance.

3228 Additional sections are informative, and add considerable information for the application
3229 developer. EXAMPLES sections provide example usage. APPLICATION USAGE sections
3230 provide additional caveats, issues, and recommendations to the developer. RATIONALE
3231 sections give additional information on the decisions made in defining the interface.

3232 FUTURE DIRECTIONS sections act as pointers to related work that may impact the interface in
3233 the future, and often cautions the developer to architect the code to account for a change in this
3234 area. Note that a future directions statement should not be taken as a commitment to adopt a
3235 feature or interface in the future.

3236 The CHANGE HISTORY section describes when the interface was introduced, and how it has
3237 changed.

3238 Option labels and margin markings in the page can be useful in guiding the application
3239 developer.

3240 **B.2 General Information**3241 **B.2.1 Use and Implementation of Functions**

3242 The information concerning the use of functions was adapted from a description in the ISO C
 3243 standard. Here is an example of how an application program can protect itself from functions
 3244 that may or may not be macros, rather than true functions:

3245 The *atoi()* function may be used in any of several ways:

- 3246 • By use of its associated header (possibly generating a macro expansion):

```
3247     #include <stdlib.h>
3248     /* ... */
3249     i = atoi(str);
```

- 3250 • By use of its associated header (assuredly generating a true function call):

```
3251     #include <stdlib.h>
3252     #undef atoi
3253     /* ... */
3254     i = atoi(str);
```

3255 or:

```
3256     #include <stdlib.h>
3257     /* ... */
3258     i = (atoi) (str);
```

- 3259 • By explicit declaration:

```
3260     extern int atoi (const char *);
3261     /* ... */
3262     i = atoi(str);
```

- 3263 • By implicit declaration:

```
3264     /* ... */
3265     i = atoi(str);
```

3266 (Assuming no function prototype is in scope. This is not allowed by the ISO C standard for
 3267 functions with variable arguments; furthermore, parameter type conversion “widening” is
 3268 subject to different rules in this case.)

3269 Note that the ISO C standard reserves names starting with ‘_’ for the compiler. Therefore, the
 3270 compiler could, for example, implement an intrinsic, built-in function *_asm_builtin_atoi()*, which
 3271 it recognized and expanded into inline assembly code. Then, in *<stdlib.h>*, there could be the
 3272 following:

```
3273     #define atoi(X) _asm_builtin_atoi(X)
```

3274 The user’s “normal” call to *atoi()* would then be expanded inline, but the implementor would
 3275 also be required to provide a callable function named *atoi()* for use when the application
 3276 requires it; for example, if its address is to be stored in a function pointer variable.

3277 **B.2.2 The Compilation Environment**3278 *B.2.2.1 POSIX.1 Symbols*

3279 This and the following section address the issue of “name space pollution”. The ISO C standard
 3280 requires that the name space beyond what it reserves not be altered except by explicit action of
 3281 the application writer. This section defines the actions to add the POSIX.1 symbols for those
 3282 headers where both the ISO C standard and POSIX.1 need to define symbols, and also where the
 3283 XSI Extension extends the base standard.

3284 When headers are used to provide symbols, there is a potential for introducing symbols that the
 3285 application writer cannot predict. Ideally, each header should only contain one set of symbols,
 3286 but this is not practical for historical reasons. Thus, the concept of feature test macros is
 3287 included. Two feature test macros are explicitly defined by IEEE Std 1003.1-200x; it is expected
 3288 that future revisions may add to this.

3289 **Note:** Feature test macros allow an application to announce to the implementation its desire to have
 3290 certain symbols and prototypes exposed. They should not be confused with the version test
 3291 macros and constants for options in `<unistd.h>` which are the implementation’s way of
 3292 announcing functionality to the application.

3293 It is further intended that these feature test macros apply only to the headers specified by
 3294 IEEE Std 1003.1-200x. Implementations are expressly permitted to make visible symbols not
 3295 specified by IEEE Std 1003.1-200x, within both POSIX.1 and other headers, under the control of
 3296 feature test macros that are not defined by IEEE Std 1003.1-200x.

3297 **The `_POSIX_C_SOURCE` Feature Test Macro**

3298 Since `_POSIX_SOURCE` specified by the POSIX.1-1990 standard did not have a value associated
 3299 with it, the `_POSIX_C_SOURCE` macro replaces it, allowing an application to inform the system
 3300 of the revision of the standard to which it conforms. This symbol will allow implementations to
 3301 support various revisions of IEEE Std 1003.1-200x simultaneously. For instance, when either
 3302 `_POSIX_SOURCE` is defined or `_POSIX_C_SOURCE` is defined as 1, the system should make
 3303 visible the same name space as permitted and required by the POSIX.1-1990 standard. When
 3304 `_POSIX_C_SOURCE` is defined, the state of `_POSIX_SOURCE` is completely irrelevant.

3305 It is expected that C bindings to future POSIX standards will define new values for
 3306 `_POSIX_C_SOURCE`, with each new value reserving the name space for that new standard, plus
 3307 all earlier POSIX standards.

3308 **The `_XOPEN_SOURCE` Feature Test Macro**

3309 The feature test macro `_XOPEN_SOURCE` is provided as the announcement mechanism for the
 3310 application that it requires functionality from the Single UNIX Specification. `_XOPEN_SOURCE`
 3311 must be defined to the value 600 before the inclusion of any header to enable the functionality in
 3312 the Single UNIX Specification. Its definition subsumes the use of `_POSIX_SOURCE` and
 3313 `_POSIX_C_SOURCE`.

3314 An extract of code from a conforming application, that appears before any `#include` statements,
 3315 is given below:

```
3316 #define _XOPEN_SOURCE 600 /* Single UNIX Specification, Version 3 */
3317 #include ...
```

3318 Note that the definition of `_XOPEN_SOURCE` with the value 600 makes the definition of
 3319 `_POSIX_C_SOURCE` redundant and it can safely be omitted.

3320 B.2.2.2 The Name Space

3321 The reservation of identifiers is paraphrased from the ISO C standard. The text is included
3322 because it needs to be part of IEEE Std 1003.1-200x, regardless of possible changes in future
3323 versions of the ISO C standard.

3324 These identifiers may be used by implementations, particularly for feature test macros.
3325 Implementations should not use feature test macro names that might be reasonably used by a
3326 standard.

3327 Including headers more than once is a reasonably common practice, and it should be carried
3328 forward from the ISO C standard. More significantly, having definitions in more than one
3329 header is explicitly permitted. Where the potential declaration is “benign” (the same definition
3330 twice) the declaration can be repeated, if that is permitted by the compiler. (This is usually true
3331 of macros, for example.) In those situations where a repetition is not benign (for example,
3332 **typedefs**), conditional compilation must be used. The situation actually occurs both within the
3333 ISO C standard and within POSIX.1: **time_t** should be in **<sys/types.h>**, and the ISO C standard
3334 mandates that it be in **<time.h>**.

3335 The area of name space pollution *versus* additions to structures is difficult because of the macro
3336 structure of C. The following discussion summarizes all the various problems with and
3337 objections to the issue.

3338 Note the phrase “user-defined macro”. Users are not permitted to define macro names (or any
3339 other name) beginning with “_**[A-Z]**”. Thus, the conflict cannot occur for symbols reserved
3340 to the vendor’s name space, and the permission to add fields automatically applies, without
3341 qualification, to those symbols.

3342 1. Data structures (and unions) need to be defined in headers by implementations to meet
3343 certain requirements of POSIX.1 and the ISO C standard.

3344 2. The structures defined by POSIX.1 are typically minimal, and any practical
3345 implementation would wish to add fields to these structures either to hold additional
3346 related information or for backwards-compatibility (or both). Future standards (and *de*
3347 *facto* standards) would also wish to add to these structures. Issues of field alignment make
3348 it impractical (at least in the general case) to simply omit fields when they are not defined
3349 by the particular standard involved.

3350 Struct **dirent** is an example of such a minimal structure (although one could argue about
3351 whether the other fields need visible names). The **st_rdev** field of most implementations’
3352 **stat** structure is a common example where extension is needed and where a conflict could
3353 occur.

3354 3. Fields in structures are in an independent name space, so the addition of such fields
3355 presents no problem to the C language itself in that such names cannot interact with
3356 identically named user symbols because access is qualified by the specific structure name.

3357 4. There is an exception to this: macro processing is done at a lexical level. Thus, symbols
3358 added to a structure might be recognized as user-provided macro names at the location
3359 where the structure is declared. This only can occur if the user-provided name is declared
3360 as a macro before the header declaring the structure is included. The user’s use of the name
3361 after the declaration cannot interfere with the structure because the symbol is hidden and
3362 only accessible through access to the structure. Presumably, the user would not declare
3363 such a macro if there was an intention to use that field name.

3364 5. Macros from the same or a related header might use the additional fields in the structure,
3365 and those field names might also collide with user macros. Although this is a less frequent
3366 occurrence, since macros are expanded at the point of use, no constraint on the order of use

3367 of names can apply.

3368 6. An “obvious” solution of using names in the reserved name space and then redefining
 3369 them as macros when they should be visible does not work because this has the effect of
 3370 exporting the symbol into the general name space. For example, given a (hypothetical)
 3371 system-provided header `<h.h>`, and two parts of a C program in `a.c` and `b.c`, in header
 3372 `<h.h>`:

```
3373     struct foo {
3374         int __i;
3375     }
3376
3377     #ifdef _FEATURE_TEST
3378     #define i __i;
3379     #endif
```

3379 In file `a.c`:

```
3380     #include h.h
3381     extern int i;
3382     ...
```

3383 In file `b.c`:

```
3384     extern int i;
3385     ...
```

3386 The symbol that the user thinks of as `i` in both files has an external name of `__i` in `a.c`; the
 3387 same symbol `i` in `b.c` has an external name `i` (ignoring any hidden manipulations the
 3388 compiler might perform on the names). This would cause a mysterious name resolution
 3389 problem when `a.o` and `b.o` are linked.

3390 Simply avoiding definition then causes alignment problems in the structure.

3391 A structure of the form:

```
3392     struct foo {
3393         union {
3394             int __i;
3395             #ifdef _FEATURE_TEST
3396             int i;
3397             #endif
3398             } __ii;
3399     }
```

3400 does not work because the name of the logical field `i` is `__ii.i`, and introduction of a macro
 3401 to restore the logical name immediately reintroduces the problem discussed previously
 3402 (although its manifestation might be more immediate because a syntax error would result
 3403 if a recursive macro did not cause it to fail first).

3404 7. A more workable solution would be to declare the structure:


```

3405     struct foo {
3406         #ifdef _FEATURE_TEST
3407             int i;
3408         #else
3409             int __i;
3410         #endif
3411     }

```

3412 However, if a macro (particularly one required by a standard) is to be defined that uses
 3413 this field, two must be defined: one that uses *i*, the other that uses *__i*. If more than one
 3414 additional field is used in a macro and they are conditional on distinct combinations of
 3415 features, the complexity goes up as 2^n .

3416 All this leaves a difficult situation: vendors must provide very complex headers to deal with
 3417 what is conceptually simple and safe—adding a field to a structure. It is the possibility of user-
 3418 provided macros with the same name that makes this difficult.

3419 Several alternatives were proposed that involved constraining the user's access to part of the
 3420 name space available to the user (as specified by the ISO C standard). In some cases, this was
 3421 only until all the headers had been included. There were two proposals discussed that failed to
 3422 achieve consensus:

- 3423 1. Limiting it for the whole program.
- 3424 2. Restricting the use of identifiers containing only uppercase letters until after all system
 3425 headers had been included. It was also pointed out that because macros might wish to
 3426 access fields of a structure (and macro expansion occurs totally at point of use) restricting
 3427 names in this way would not protect the macro expansion, and thus the solution was
 3428 inadequate.

3429 It was finally decided that reservation of symbols would occur, but as constrained.

3430 The current wording also allows the addition of fields to a structure, but requires that user
 3431 macros of the same name not interfere. This allows vendors to do one of the following:

- 3432 • Not create the situation (do not extend the structures with user-accessible names or use the
 3433 solution in (7) above)
- 3434 • Extend their compilers to allow some way of adding names to structures and macros safely

3435 There are at least two ways that the compiler might be extended: add new preprocessor
 3436 directives that turn off and on macro expansion for certain symbols (without changing the value
 3437 of the macro) and a function or lexical operation that suppresses expansion of a word. The latter
 3438 seems more flexible, particularly because it addresses the problem in macros as well as in
 3439 declarations.

3440 The following seems to be a possible implementation extension to the C language that will do
 3441 this: any token that during macro expansion is found to be preceded by three '#' symbols shall
 3442 not be further expanded in exactly the same way as described for macros that expand to their
 3443 own name as in Section 3.8.3.4 of the ISO C standard. A vendor may also wish to implement this
 3444 as an operation that is lexically a function, which might be implemented as:

```

3445     #define __safe_name(x) ###x

```

3446 Using a function notation would insulate vendors from changes in standards until such a
 3447 functionality is standardized (if ever). Standardization of such a function would be valuable
 3448 because it would then permit third parties to take advantage of it portably in software they may
 3449 supply.

3450 The symbols that are “explicitly permitted, but not required by IEEE Std 1003.1-200x” include
 3451 those classified below. (That is, the symbols classified below might, but are not required to, be
 3452 present when `_POSIX_C_SOURCE` is defined to have the value 20010xL.)

- 3453 • Symbols in `<limits.h>` and `<unistd.h>` that are defined to indicate support for options or
 3454 limits that are constant at compile-time.
- 3455 • Symbols in the name space reserved for the implementation by the ISO C standard.
- 3456 • Symbols in a name space reserved for a particular type of extension (for example, type names
 3457 ending with `_t` in `<sys/types.h>`).
- 3458 • Additional members of structures or unions whose names do not reduce the name space
 3459 reserved for applications.

3460 Since both implementations and future revisions of IEEE Std 1003.1-200x and other POSIX
 3461 standards may use symbols in the reserved spaces described in these tables, there is a potential
 3462 for name space clashes. To avoid future name space clashes when adding symbols,
 3463 implementations should not use the `posix_`, `POSIX_`, or `_POSIX_` prefixes.

3464 B.2.3 Error Numbers

3465 It was the consensus of the standard developers that to allow the conformance document to
 3466 state that an error occurs and under what conditions, but to disallow a statement that it never
 3467 occurs, does not make sense. It could be implied by the current wording that this is allowed, but
 3468 to reduce the possibility of future interpretation requests, it is better to make an explicit
 3469 statement.

3470 The ISO C standard requires that `errno` be an assignable lvalue. Originally, the definition in |
 3471 POSIX.1 was stricter than that in the ISO C standard, `extern int errno`, in order to support |
 3472 historical usage. In a multi-threaded environment, implementing `errno` as a global variable
 3473 results in non-deterministic results when accessed. It is required, however, that `errno` work as a
 3474 per-thread error reporting mechanism. In order to do this, a separate `errno` value has to be
 3475 maintained for each thread. The following section discusses the various alternative solutions
 3476 that were considered.

3477 In order to avoid this problem altogether for new functions, these functions avoid using `errno`
 3478 and, instead, return the error number directly as the function return value; a return value of zero
 3479 indicates that no error was detected.

3480 For any function that can return errors, the function return value is not used for any purpose
 3481 other than for reporting errors. Even when the output of the function is scalar, it is passed
 3482 through a function argument. While it might have been possible to allow some scalar outputs to
 3483 be coded as negative function return values and mixed in with positive error status returns, this
 3484 was rejected—using the return value for a mixed purpose was judged to be of limited use and
 3485 error prone.

3486 Checking the value of `errno` alone is not sufficient to determine the existence or type of an error,
 3487 since it is not required that a successful function call clear `errno`. The variable `errno` should only
 3488 be examined when the return value of a function indicates that the value of `errno` is meaningful.
 3489 In that case, the function is required to set the variable to something other than zero.

3490 The variable `errno` shall never be set to zero by any function call; to do so would contradict the
 3491 ISO C standard.

3492 POSIX.1 requires (in the ERRORS sections of function descriptions) certain error values to be set
 3493 in certain conditions because many existing applications depend on them. Some error numbers,
 3494 such as [EFAULT], are entirely implementation-defined and are noted as such in their

3495 description in the ERRORS section. This section otherwise allows wide latitude to the
 3496 implementation in handling error reporting.

3497 Some of the ERRORS sections in IEEE Std 1003.1-200x have two subsections. The first:

3498 “The function shall fail if:”

3499 could be called the “mandatory” section.

3500 The second:

3501 “The function may fail if:”

3502 could be informally known as the “optional” section.

3503 Attempting to infer the quality of an implementation based on whether it detects optional error
 3504 conditions is not useful.

3505 Following each one-word symbolic name for an error, there is a description of the error. The
 3506 rationale for some of the symbolic names follows:

3507 [ECANCELED] This spelling was chosen as being more common.

3508 [EFAULT] Most historical implementations do not catch an error and set *errno* when an
 3509 invalid address is given to the functions *wait()*, *time()*, or *times()*. Some
 3510 implementations cannot reliably detect an invalid address. And most systems
 3511 that detect invalid addresses will do so only for a system call, not for a library
 3512 routine.

3513 [EFTYPE] This error code was proposed in earlier proposals as “Inappropriate operation
 3514 for file type”, meaning that the operation requested is not appropriate for the
 3515 file specified in the function call. This code was proposed, although the same
 3516 idea was covered by [ENOTTY], because the connotations of the name would
 3517 be misleading. It was pointed out that the *fcntl()* function uses the error code
 3518 [EINVAL] for this notion, and hence all instances of [EFTYPE] were changed
 3519 to this code.

3520 [EINTR] POSIX.1 prohibits conforming implementations from restarting interrupted
 3521 system calls. However, it does not require that [EINTR] be returned when
 3522 another legitimate value may be substituted; for example, a partial transfer
 3523 count when *read()* or *write()* are interrupted. This is only given when the
 3524 signal catching function returns normally as opposed to returns by
 3525 mechanisms like *longjmp()* or *siglongjmp()*.

3526 [ELOOP] In specifying conditions under which implementations would generate this
 3527 error, the following goals were considered:

- 3528 • To ensure that actual loops are detected, including loops that result from
 3529 symbolic links across distributed file systems.
- 3530 • To ensure that during pathname resolution an application can rely on the
 3531 ability to follow at least {SYMLOOP_MAX} symbolic links in the absence
 3532 of a loop.
- 3533 • To allow implementations to provide the capability of traversing more
 3534 than {SYMLOOP_MAX} symbolic links in the absence of a loop.
- 3535 • To allow implementations to detect loops and generate the error prior to
 3536 encountering {SYMLOOP_MAX} symbolic links.

3537	[ENAMETOOLONG]	
3538		When a symbolic link is encountered during pathname resolution, the
3539		contents of that symbolic link are used to create a new pathname. The
3540		standard developers intended to allow, but not require, that implementations
3541		enforce the restriction of {PATH_MAX} on the result of this pathname
3542		substitution.
3543	[ENOMEM]	The term <i>main memory</i> is not used in POSIX.1 because it is implementation-
3544		defined.
3545	[ENOTSUP]	This error code is to be used when an implementation chooses to implement
3546		the required functionality of IEEE Std 1003.1-200x but does not support
3547		optional facilities defined by IEEE Std 1003.1-200x. The return of [ENOSYS] is
3548		to be taken to indicate that the function of the interface is not supported at all;
3549		the function will always fail with this error code.
3550	[ENOTTY]	The symbolic name for this error is derived from a time when device control
3551		was done by <i>ioctl()</i> and that operation was only permitted on a terminal
3552		interface. The term <i>TTY</i> is derived from <i>teletypewriter</i> , the devices to which
3553		this error originally applied.
3554	[EOVERFLOW]	Most of the uses of this error code are related to large file support. Typically,
3555		these cases occur on systems which support multiple programming
3556		environments with different sizes for <i>off_t</i> , but they may also occur in
3557		connection with remote file systems.
3558		In addition, when different programming environments have different widths
3559		for types such as <i>int</i> and <i>uid_t</i> , several functions may encounter a condition
3560		where a value in a particular environment is too wide to be represented. In
3561		that case, this error should be raised. For example, suppose the currently
3562		running process has 64-bit <i>int</i> , and file descriptor 9223372036854775807 is
3563		open and does not have the close-on-exec flag set. If the process then uses
3564		<i>execl()</i> to <i>exec</i> a file compiled in a programming environment with 32-bit <i>int</i> ,
3565		the call to <i>execl()</i> can fail with <i>errno</i> set to [EOVERFLOW]. A similar failure
3566		can occur with <i>execl()</i> if any of the user IDs or any of the group IDs to be
3567		assigned to the new process image are out of range for the executed file's
3568		programming environment.
3569		Note, however, that this condition cannot occur for functions that are
3570		explicitly described as always being successful, such as <i>getpid()</i> .
3571	[EPIPE]	This condition normally generates the signal SIGPIPE; the error is returned if
3572		the signal does not terminate the process.
3573	[EROFS]	In historical implementations, attempting to <i>unlink()</i> or <i>rmdir()</i> a mount point
3574		would generate an [EBUSY] error. An implementation could be envisioned
3575		where such an operation could be performed without error. In this case, if
3576		<i>either</i> the directory entry or the actual data structures reside on a read-only file
3577		system, [EROFS] is the appropriate error to generate. (For example, changing
3578		the link count of a file on a read-only file system could not be done, as is
3579		required by <i>unlink()</i> , and thus an error should be reported.)
3580		Three error numbers, [EDOM], [EILSEQ], and [ERANGE], were added to this section primarily
3581		for consistency with the ISO C standard.

3582 **Alternative Solutions for Per-Thread *errno***

3583 The usual implementation of *errno* as a single global variable does not work in a multi-threaded
 3584 environment. In such an environment, a thread may make a POSIX.1 call and get a -1 error
 3585 return, but before that thread can check the value of *errno*, another thread might have made a
 3586 second POSIX.1 call that also set *errno*. This behavior is unacceptable in robust programs. There
 3587 were a number of alternatives that were considered for handling the *errno* problem:

- 3588 • Implement *errno* as a per-thread integer variable.
- 3589 • Implement *errno* as a service that can access the per-thread error number.
- 3590 • Change all POSIX.1 calls to accept an extra status argument and avoid setting *errno*.
- 3591 • Change all POSIX.1 calls to raise a language exception.

3592 The first option offers the highest level of compatibility with existing practice but requires
 3593 special support in the linker, compiler, and/or virtual memory system to support the new
 3594 concept of thread private variables. When compared with current practice, the third and fourth
 3595 options are much cleaner, more efficient, and encourage a more robust programming style, but
 3596 they require new versions of all of the POSIX.1 functions that might detect an error. The second
 3597 option offers compatibility with existing code that uses the `<errno.h>` header to define the
 3598 symbol *errno*. In this option, *errno* may be a macro defined:

```
3599     #define errno  (*__errno())
3600     extern int    *__errno();
```

3601 This option may be implemented as a per-thread variable whereby an *errno* field is allocated in
 3602 the user space object representing a thread, and whereby the function `__errno()` makes a system
 3603 call to determine the location of its user space object and returns the address of the *errno* field of
 3604 that object. Another implementation, one that avoids calling the kernel, involves allocating
 3605 stacks in chunks. The stack allocator keeps a side table indexed by chunk number containing a
 3606 pointer to the thread object that uses that chunk. The `__errno()` function then looks at the stack
 3607 pointer, determines the chunk number, and uses that as an index into the chunk table to find its
 3608 thread object and thus its private value of *errno*. On most architectures, this can be done in four
 3609 to five instructions. Some compilers may wish to implement `__errno()` inline to improve
 3610 performance.

3611 **Disallowing Return of the [EINTR] Error Code**

3612 Many blocking interfaces defined by IEEE Std 1003.1-200x may return [EINTR] if interrupted
 3613 during their execution by a signal handler. Blocking interfaces introduced under the Threads
 3614 option do not have this property. Instead, they require that the interface appear to be atomic
 3615 with respect to interruption. In particular, clients of blocking interfaces need not handle any
 3616 possible [EINTR] return as a special case since it will never occur. If it is necessary to restart
 3617 operations or complete incomplete operations following the execution of a signal handler, this is
 3618 handled by the implementation, rather than by the application.

3619 Requiring applications to handle [EINTR] errors on blocking interfaces has been shown to be a
 3620 frequent source of often unreproducible bugs, and it adds no compelling value to the available
 3621 functionality. Thus, blocking interfaces introduced for use by multi-threaded programs do not
 3622 use this paradigm. In particular, in none of the functions `flockfile()`, `pthread_cond_timedwait()`,
 3623 `pthread_cond_wait()`, `pthread_join()`, `pthread_mutex_lock()`, and `sigwait()` did providing [EINTR]
 3624 returns add value, or even particularly make sense. Thus, these functions do not provide for an
 3625 [EINTR] return, even when interrupted by a signal handler. The same arguments can be applied
 3626 to `sem_wait()`, `sem_trywait()`, `sigwaitinfo()`, and `sigtimedwait()`, but implementations are
 3627 permitted to return [EINTR] error codes for these functions for compatibility with earlier

3628 versions of IEEE Std 1003.1-200x. Applications cannot rely on calls to these functions returning
3629 [EINTR] error codes when signals are delivered to the calling thread, but they should allow for
3630 the possibility.

3631 *B.2.3.1 Additional Error Numbers*

3632 The ISO C standard defines the name space for implementations to add additional error
3633 numbers.

3634 **B.2.4 Signal Concepts**

3635 Historical implementations of signals, using the *signal()* function, have shortcomings that make
3636 them unreliable for many application uses. Because of this, a new signal mechanism, based very
3637 closely on the one of 4.2 BSD and 4.3 BSD, was added to POSIX.1.

3638 **Signal Names**

3639 The restriction on the actual type used for **sigset_t** is intended to guarantee that these objects can
3640 always be assigned, have their address taken, and be passed as parameters by value. It is not
3641 intended that this type be a structure including pointers to other data structures, as that could
3642 impact the portability of applications performing such operations. A reasonable implementation
3643 could be a structure containing an array of some integer type.

3644 The signals described in IEEE Std 1003.1-200x must have unique values so that they may be
3645 named as parameters of **case** statements in the body of a C language **switch** clause. However,
3646 implementation-defined signals may have values that overlap with each other or with signals
3647 specified in IEEE Std 1003.1-200x. An example of this is SIGABRT, which traditionally overlaps
3648 some other signal, such as SIGIOT.

3649 SIGKILL, SIGTERM, SIGUSR1, and SIGUSR2 are ordinarily generated only through the explicit
3650 use of the *kill()* function, although some implementations generate SIGKILL under
3651 extraordinary circumstances. SIGTERM is traditionally the default signal sent by the *kill*
3652 command.

3653 The signals SIGBUS, SIGEMT, SIGIOT, SIGTRAP, and SIGSYS were omitted from POSIX.1
3654 because their behavior is implementation-defined and could not be adequately categorized.
3655 Conforming implementations may deliver these signals, but must document the circumstances
3656 under which they are delivered and note any restrictions concerning their delivery. The signals
3657 SIGFPE, SIGILL, and SIGSEGV are similar in that they also generally result only from
3658 programming errors. They were included in POSIX.1 because they do indicate three relatively
3659 well-categorized conditions. They are all defined by the ISO C standard and thus would have to
3660 be defined by any system with a ISO C standard binding, even if not explicitly included in
3661 POSIX.1.

3662 There is very little that a Conforming POSIX.1 Application can do by catching, ignoring, or
3663 masking any of the signals SIGILL, SIGTRAP, SIGIOT, SIGEMT, SIGBUS, SIGSEGV, SIGSYS, or
3664 SIGFPE. They will generally be generated by the system only in cases of programming errors.
3665 While it may be desirable for some robust code (for example, a library routine) to be able to
3666 detect and recover from programming errors in other code, these signals are not nearly sufficient
3667 for that purpose. One portable use that does exist for these signals is that a command interpreter
3668 can recognize them as the cause of a process' termination (with *wait()*) and print an appropriate
3669 message. The mnemonic tags for these signals are derived from their PDP-11 origin.

3670 The signals SIGSTOP, SIGTSTP, SIGTTIN, SIGTTOU, and SIGCONT are provided for job control
3671 and are unchanged from 4.2 BSD. The signal SIGCHLD is also typically used by job control
3672 shells to detect children that have terminated or, as in 4.2 BSD, stopped.

3673 Some implementations, including System V, have a signal named SIGCLD, which is similar to
3674 SIGCHLD in 4.2 BSD. POSIX.1 permits implementations to have a single signal with both
3675 names. POSIX.1 carefully specifies ways in which conforming applications can avoid the
3676 semantic differences between the two different implementations. The name SIGCHLD was
3677 chosen for POSIX.1 because most current application usages of it can remain unchanged in
3678 conforming applications. SIGCLD in System V has more cases of semantics that POSIX.1 does
3679 not specify, and thus applications using it are more likely to require changes in addition to the
3680 name change.

3681 The signals SIGUSR1 and SIGUSR2 are commonly used by applications for notification of
3682 exceptional behavior and are described as “reserved as application-defined” so that such use is
3683 not prohibited. Implementations should not generate SIGUSR1 or SIGUSR2, except when
3684 explicitly requested by *kill()*. It is recommended that libraries not use these two signals, as such
3685 use in libraries could interfere with their use by applications calling the libraries. If such use is
3686 unavoidable, it should be documented. It is prudent for non-portable libraries to use non-
3687 standard signals to avoid conflicts with use of standard signals by portable libraries.

3688 There is no portable way for an application to catch or ignore non-standard signals. Some
3689 implementations define the range of signal numbers, so applications can install signal-catching
3690 functions for all of them. Unfortunately, implementation-defined signals often cause problems
3691 when caught or ignored by applications that do not understand the reason for the signal. While
3692 the desire exists for an application to be more robust by handling all possible signals (even those
3693 only generated by *kill()*), no existing mechanism was found to be sufficiently portable to include
3694 in POSIX.1. The value of such a mechanism, if included, would be diminished given that
3695 SIGKILL would still not be catchable.

3696 A number of new signal numbers are reserved for applications because the two user signals
3697 defined by POSIX.1 are insufficient for many realtime applications. A range of signal numbers is
3698 specified, rather than an enumeration of additional reserved signal names, because different
3699 applications and application profiles will require a different number of application signals. It is
3700 not desirable to burden all application domains and therefore all implementations with the
3701 maximum number of signals required by all possible applications. Note that in this context,
3702 signal numbers are essentially different signal priorities.

3703 The relatively small number of required additional signals, `{_POSIX_RTSIG_MAX}`, was chosen
3704 so as not to require an unreasonably large signal mask/set. While this number of signals defined
3705 in POSIX.1 will fit in a single 32-bit word signal mask, it is recognized that most existing
3706 implementations define many more signals than are specified in POSIX.1 and, in fact, many
3707 implementations have already exceeded 32 signals (including the “null signal”). Support of
3708 `{_POSIX_RTSIG_MAX}` additional signals may push some implementation over the single 32-bit
3709 word line, but is unlikely to push any implementations that are already over that line beyond the
3710 64-signal line.

3711 B.2.4.1 Signal Generation and Delivery

3712 The terms defined in this section are not used consistently in documentation of historical
3713 systems. Each signal can be considered to have a lifetime beginning with *generation* and ending
3714 with *delivery* or *acceptance*. The POSIX.1 definition of *delivery* does not exclude ignored signals;
3715 this is considered a more consistent definition. This revised text in several parts of
3716 IEEE Std 1003.1-200x clarifies the distinct semantics of asynchronous signal *delivery* and
3717 synchronous signal *acceptance*. The previous wording attempted to categorize both under the
3718 term *delivery*, which led to conflicts over whether the effects of asynchronous signal delivery
3719 applied to synchronous signal acceptance.

3720 Signals generated for a process are delivered to only one thread. Thus, if more than one thread is
3721 eligible to receive a signal, one has to be chosen. The choice of threads is left entirely up to the
3722 implementation both to allow the widest possible range of conforming implementations and to
3723 give implementations the freedom to deliver the signal to the “easiest possible” thread should
3724 there be differences in ease of delivery between different threads.

3725 Note that should multiple delivery among cooperating threads be required by an application,
3726 this can be trivially constructed out of the provided single-delivery semantics. The construction
3727 of a *sigwait_multiple()* function that accomplishes this goal is presented with the rationale for
3728 *sigwaitinfo()*.

3729 Implementations should deliver unblocked signals as soon after they are generated as possible.
3730 However, it is difficult for POSIX.1 to make specific requirements about this, beyond those in
3731 *kill()* and *sigprocmask()*. Even on systems with prompt delivery, scheduling of higher priority
3732 processes is always likely to cause delays.

3733 In general, the interval between the generation and delivery of unblocked signals cannot be
3734 detected by an application. Thus, references to pending signals generally apply to blocked,
3735 pending signals. An implementation registers a signal as pending on the process when no thread
3736 has the signal unblocked and there are no threads blocked in a *sigwait()* function for that signal.
3737 Thereafter, the implementation delivers the signal to the first thread that unblocks the signal or
3738 calls a *sigwait()* function on a signal set containing this signal rather than choosing the recipient
3739 thread at the time the signal is sent.

3740 In the 4.3 BSD system, signals that are blocked and set to SIG_IGN are discarded immediately
3741 upon generation. For a signal that is ignored as its default action, if the action is SIG_DFL and
3742 the signal is blocked, a generated signal remains pending. In the 4.1 BSD system and in
3743 System V, Release 3, two other implementations that support a somewhat similar signal
3744 mechanism, all ignored, blocked signals remain pending if generated. Because it is not normally
3745 useful for an application to simultaneously ignore and block the same signal, it was unnecessary
3746 for POSIX.1 to specify behavior that would invalidate any of the historical implementations.

3747 There is one case in some historical implementations where an unblocked, pending signal does
3748 not remain pending until it is delivered. In the System V implementation of *signal()*, pending
3749 signals are discarded when the action is set to SIG_DFL or a signal-catching routine (as well as to
3750 SIG_IGN). Except in the case of setting SIGCHLD to SIG_DFL, implementations that do this do
3751 not conform completely to POSIX.1. Some earlier proposals for POSIX.1 explicitly stated this,
3752 but these statements were redundant due to the requirement that functions defined by POSIX.1
3753 not change attributes of processes defined by POSIX.1 except as explicitly stated.

3754 POSIX.1 specifically states that the order in which multiple, simultaneously pending signals are
3755 delivered is unspecified. This order has not been explicitly specified in historical
3756 implementations, but has remained quite consistent and been known to those familiar with the
3757 implementations. Thus, there have been cases where applications (usually system utilities) have
3758 been written with explicit or implicit dependencies on this order. Implementors and others
3759 porting existing applications may need to be aware of such dependencies.

3760 When there are multiple pending signals that are not blocked, implementations should arrange
3761 for the delivery of all signals at once, if possible. Some implementations stack calls to all pending
3762 signal-catching routines, making it appear that each signal-catcher was interrupted by the next
3763 signal. In this case, the implementation should ensure that this stacking of signals does not
3764 violate the semantics of the signal masks established by *sigaction()*. Other implementations
3765 process at most one signal when the operating system is entered, with remaining signals saved
3766 for later delivery. Although this practice is widespread, this behavior is neither standardized
3767 nor endorsed. In either case, implementations should attempt to deliver signals associated with
3768 the current state of the process (for example, SIGFPE) before other signals, if possible.

3769 In 4.2 BSD and 4.3 BSD, it is not permissible to ignore or explicitly block SIGCONT, because if
 3770 blocking or ignoring this signal prevented it from continuing a stopped process, such a process
 3771 could never be continued (only killed by SIGKILL). However, 4.2 BSD and 4.3 BSD do block
 3772 SIGCONT during execution of its signal-catching function when it is caught, creating exactly
 3773 this problem. A proposal was considered to disallow catching SIGCONT in addition to ignoring
 3774 and blocking it, but this limitation led to objections. The consensus was to require that
 3775 SIGCONT always continue a stopped process when generated. This removed the need to
 3776 disallow ignoring or explicit blocking of the signal; note that SIG_IGN and SIG_DFL are
 3777 equivalent for SIGCONT .

3778 B.2.4.2 Realtime Signal Generation and Delivery

3779 The Realtime Signals Extension option to POSIX.1 signal generation and delivery behavior is
 3780 required for the following reasons:

- 3781 • The **sigevent** structure is used by other POSIX.1 functions that result in asynchronous event
 3782 notifications to specify the notification mechanism to use and other information needed by
 3783 the notification mechanism. IEEE Std 1003.1-200x defines only three symbolic values for the
 3784 notification mechanism. SIGEV_NONE is used to indicate that no notification is required
 3785 when the event occurs. This is useful for applications that use asynchronous I/O with polling
 3786 for completion. SIGEV_SIGNAL indicates that a signal shall be generated when the event
 3787 occurs. SIGEV_NOTIFY provides for “callback functions” for asynchronous notifications
 3788 done by a function call within the context of a new thread. This provides a multi-threaded
 3789 process a more natural means of notification than signals. The primary difficulty with
 3790 previous notification approaches has been to specify the environment of the notification
 3791 routine.
 - 3792 — One approach is to limit the notification routine to call only functions permitted in a
 3793 signal handler. While the list of permissible functions is clearly stated, this is overly
 3794 restrictive.
 - 3795 — A second approach is to define a new list of functions or classes of functions that are
 3796 explicitly permitted or not permitted. This would give a programmer more lists to deal
 3797 with, which would be awkward.
 - 3798 — The third approach is to define completely the environment for execution of the
 3799 notification function. A clear definition of an execution environment for notification is
 3800 provided by executing the notification function in the environment of a newly created
 3801 thread.

3802 Implementations may support additional notification mechanisms by defining new values
 3803 for *sigev_notify*.

3804 For a notification type of SIGEV_SIGNAL, the other members of the **sigevent** structure
 3805 defined by IEEE Std 1003.1-200x specify the realtime signal—that is, the signal number and
 3806 application-defined value that differentiates between occurrences of signals with the same
 3807 number—that will be generated when the event occurs. The structure is defined in
 3808 <signal.h>, even though the structure is not directly used by any of the signal functions,
 3809 because it is part of the signals interface used by the POSIX.1b “client functions”. When the
 3810 client functions include <signal.h> to define the signal names, the **sigevent** structure will
 3811 also be defined.

3812 An application-defined value passed to the signal handler is used to differentiate between
 3813 different “events” instead of requiring that the application use different signal numbers for
 3814 several reasons:

- 3815 — Realtime applications potentially handle a very large number of different events.
 3816 Requiring that implementations support a correspondingly large number of distinct
 3817 signal numbers will adversely impact the performance of signal delivery because the
 3818 signal masks to be manipulated on entry and exit to the handlers will become large.
- 3819 — Event notifications are prioritized by signal number (the rationale for this is explained in
 3820 the following paragraphs) and the use of different signal numbers to differentiate
 3821 between the different event notifications overloads the signal number more than has
 3822 already been done. It also requires that the application writer make arbitrary assignments
 3823 of priority to events that are logically of equal priority.
- 3824 A union is defined for the application-defined value so that either an integer constant or a
 3825 pointer can be portably passed to the signal-catching function. On some architectures a
 3826 pointer cannot be cast to an `int` and *vice versa*.
- 3827 Use of a structure here with an explicit notification type discriminant rather than explicit
 3828 parameters to realtime functions, or embedded in other realtime structures, provides for
 3829 future extensions to IEEE Std 1003.1-200x. Additional, perhaps more efficient, notification
 3830 mechanisms can be supported for existing realtime function interfaces, such as timers and
 3831 asynchronous I/O, by extending the `sigevent` structure appropriately. The existing realtime
 3832 function interfaces will not have to be modified to use any such new notification mechanism.
 3833 The revised text concerning the `SIGEV_SIGNAL` value makes consistent the semantics of the
 3834 members of the `sigevent` structure, particularly in the definitions of `lio_listio()` and
 3835 `aio_fsync()`. For uniformity, other revisions cause this specification to be referred to rather
 3836 than inaccurately duplicated in the descriptions of functions and structures using the
 3837 `sigevent` structure. The revised wording does not relax the requirement that the signal
 3838 number be in the range `SIGRTMIN` to `SIGRTMAX` to guarantee queuing and passing of the
 3839 application value, since that requirement is still implied by the signal names.
- 3840 • IEEE Std 1003.1-200x is intentionally vague on whether “non-realtime” signal-generating
 3841 mechanisms can result in a `siginfo_t` being supplied to the handler on delivery. In one
 3842 existing implementation, a `siginfo_t` is posted on signal generation, even though the
 3843 implementation does not support queuing of multiple occurrences of a signal. It is not the
 3844 intent of IEEE Std 1003.1-200x to preclude this, independent of the mandate to define signals
 3845 that do support queuing. Any interpretation that appears to preclude this is a mistake in the
 3846 reading or writing of the standard.
 - 3847 • Signals handled by realtime signal handlers might be generated by functions or conditions
 3848 that do not allow the specification of an application-defined value and do not queue.
 3849 IEEE Std 1003.1-200x specifies the `si_code` member of the `siginfo_t` structure used in existing
 3850 practice and defines additional codes so that applications can detect whether an application-
 3851 defined value is present or not. The code `SI_USER` for `kill()`-generated signals is adopted
 3852 from existing practice.
 - 3853 • The `sigaction()` `sa_flags` value `SA_SIGINFO` tells the implementation that the signal-catching
 3854 function expects two additional arguments. When the flag is not set, a single argument, the
 3855 signal number, is passed as specified by IEEE Std 1003.1-200x. Although IEEE Std 1003.1-200x
 3856 does not explicitly allow the `info` argument to the handler function to be `NULL`, this is
 3857 existing practice. This provides for compatibility with programs whose signal-catching
 3858 functions are not prepared to accept the additional arguments. IEEE Std 1003.1-200x is
 3859 explicitly unspecified as to whether signals actually queue when `SA_SIGINFO` is not set for a
 3860 signal, as there appear to be no benefits to applications in specifying one behavior or another.
 3861 One existing implementation queues a `siginfo_t` on each signal generation, unless the signal
 3862 is already pending, in which case the implementation discards the new `siginfo_t`; that is, the
 3863 queue length is never greater than one. This implementation only examines `SA_SIGINFO` on

3864 signal delivery, discarding the queued `siginfo_t` if its delivery was not requested.

3865 IEEE Std 1003.1-200x specifies several new values for the `si_code` member of the `siginfo_t`
3866 structure. In existing practice, a `si_code` value of less than or equal to zero indicates that the
3867 signal was generated by a process via the `kill()` function. In existing practice, values of `si_code`
3868 that provide additional information for implementation-generated signals, such as SIGFPE or
3869 SIGSEGV, are all positive. Thus, if implementations define the new constants specified in
3870 IEEE Std 1003.1-200x to be negative numbers, programs written to use existing practice will
3871 not break. IEEE Std 1003.1-200x chose not to attempt to specify existing practice values of
3872 `si_code` other than SI_USER both because it was deemed beyond the scope of
3873 IEEE Std 1003.1-200x and because many of the values in existing practice appear to be
3874 platform and implementation-defined. But, IEEE Std 1003.1-200x does specify that if an
3875 implementation—for example, one that does not have existing practice in this area—chooses
3876 to define additional values for `si_code`, these values have to be different from the values of the
3877 symbols specified by IEEE Std 1003.1-200x. This will allow conforming applications to
3878 differentiate between signals generated by one of the POSIX.1b asynchronous events and
3879 those generated by other implementation events in a manner compatible with existing
3880 practice.

3881 The unique values of `si_code` for the POSIX.1b asynchronous events have implications for
3882 implementations of, for example, asynchronous I/O or message passing in user space library
3883 code. Such an implementation will be required to provide a hidden interface to the signal
3884 generation mechanism that allows the library to specify the standard values of `si_code`.

3885 Existing practice also defines additional members of `siginfo_t`, such as the process ID and
3886 user ID of the sending process for `kill()`-generated signals. These members were deemed not
3887 necessary to meet the requirements of realtime applications and are not specified by
3888 IEEE Std 1003.1-200x. Neither are they precluded.

3889 The third argument to the signal-catching function, `context`, is left undefined by
3890 IEEE Std 1003.1-200x, but is specified in the interface because it matches existing practice for
3891 the SA_SIGINFO flag. It was considered undesirable to require a separate implementation
3892 for SA_SIGINFO for POSIX conformance on implementations that already support the two
3893 additional parameters.

- 3894 • The requirement to deliver lower numbered signals in the range SIGRTMIN to SIGRTMAX
3895 first, when multiple unblocked signals are pending, results from several considerations:
 - 3896 — A method is required to prioritize event notifications. The signal number was chosen
3897 instead of, for instance, associating a separate priority with each request, because an
3898 implementation has to check pending signals at various points and select one for delivery
3899 when more than one is pending. Specifying a selection order is the minimal additional
3900 semantic that will achieve prioritized delivery. If a separate priority were to be associated
3901 with queued signals, it would be necessary for an implementation to search all non-
3902 empty, non-blocked signal queues and select from among them the pending signal with
3903 the highest priority. This would significantly increase the cost of and decrease the
3904 determinism of signal delivery.
 - 3905 — Given the specified selection of the lowest numeric unblocked pending signal,
3906 preemptive priority signal delivery can be achieved using signal numbers and signal
3907 masks by ensuring that the `sa_mask` for each signal number blocks all signals with a
3908 higher numeric value.

3909 For realtime applications that want to use only the newly defined realtime signal numbers
3910 without interference from the standard signals, this can be achieved by blocking all of the
3911 standard signals in the process signal mask and in the `sa_mask` installed by the signal

3912 action for the realtime signal handlers.

3913 IEEE Std 1003.1-200x explicitly leaves unspecified the ordering of signals outside of the range
3914 of realtime signals and the ordering of signals within this range with respect to those outside
3915 the range. It was believed that this would unduly constrain implementations or standards in
3916 the future definition of new signals.

3917 B.2.4.3 Signal Actions

3918 Early proposals mentioned SIGCONT as a second exception to the rule that signals are not
3919 delivered to stopped processes until continued. Because IEEE Std 1003.1-200x now specifies that
3920 SIGCONT causes the stopped process to continue when it is generated, delivery of SIGCONT is
3921 not prevented because a process is stopped, even without an explicit exception to this rule.

3922 Ignoring a signal by setting the action to SIG_IGN (or SIG_DFL for signals whose default action
3923 is to ignore) is not the same as installing a signal-catching function that simply returns. Invoking
3924 such a function will interrupt certain system functions that block processes (for example, *wait()*,
3925 *sigsuspend()*, *pause()*, *read()*, *write()*) while ignoring a signal has no such effect on the process.

3926 Historical implementations discard pending signals when the action is set to SIG_IGN.
3927 However, they do not always do the same when the action is set to SIG_DFL and the default
3928 action is to ignore the signal. IEEE Std 1003.1-200x requires this for the sake of consistency and
3929 also for completeness, since the only signal this applies to is SIGCHLD, and
3930 IEEE Std 1003.1-200x disallows setting its action to SIG_IGN.

3931 Some implementations (System V, for example) assign different semantics for SIGCLD |
3932 depending on whether the action is set to SIG_IGN or SIG_DFL. Since POSIX.1 requires that the |
3933 default action for SIGCHLD be to ignore the signal, applications should always set the action to |
3934 SIG_DFL in order to avoid SIGCHLD.

3935 Whether or not an implementation allows SIG_IGN as a SIGCHLD disposition to be inherited |
3936 across a call to one of the *exec* family of functions or *posix_spawn()* is explicitly left as |
3937 unspecified. This change was made as a result of IEEE PASC Interpretation 1003.1 #132, and |
3938 permits the implementation to decide between the following alternatives: |

- 3939 • Unconditionally leave SIGCHLD set to SIG_IGN, in which case the implementation would |
3940 not allow applications that assume inheritance of SIG_DFL to conform to |
3941 IEEE Std 1003.1-200x without change. The implementation would, however, retain an ability |
3942 to control applications that create child processes but never call on the *wait* family of |
3943 functions, potentially filling up the process table. |
- 3944 • Unconditionally reset SIGCHLD to SIG_DFL, in which case the implementation would allow |
3945 applications that assume inheritance of SIG_DFL to conform. The implementation would, |
3946 however, lose an ability to control applications that spawn child processes but never reap |
3947 them. |
- 3948 • Provide some mechanism, not specified in IEEE Std 1003.1-200x, to control inherited |
3949 SIGCHLD dispositions. |

3950 Some implementations (System V, for example) will deliver a SIGCLD signal immediately when |
3951 a process establishes a signal-catching function for SIGCLD when that process has a child that |
3952 has already terminated. Other implementations, such as 4.3 BSD, do not generate a new |
3953 SIGCHLD signal in this way. In general, a process should not attempt to alter the signal action |
3954 for the SIGCHLD signal while it has any outstanding children. However, it is not always |
3955 possible for a process to avoid this; for example, shells sometimes start up processes in pipelines |
3956 with other processes from the pipeline as children. Processes that cannot ensure that they have |
3957 no children when altering the signal action for SIGCHLD thus need to be prepared for, but not

3958 depend on, generation of an immediate SIGCHLD signal.

3959 The default action of the stop signals (SIGSTOP , SIGTSTP, SIGTTIN, SIGTTOU) is to stop a
 3960 process that is executing. If a stop signal is delivered to a process that is already stopped, it has
 3961 no effect. In fact, if a stop signal is generated for a stopped process whose signal mask blocks the
 3962 signal, the signal will never be delivered to the process since the process must receive a
 3963 SIGCONT, which discards all pending stop signals, in order to continue executing.

3964 The SIGCONT signal shall continue a stopped process even if SIGCONT is blocked (or ignored).
 3965 However, if a signal-catching routine has been established for SIGCONT, it will not be entered
 3966 until SIGCONT is unblocked.

3967 If a process in an orphaned process group stops, it is no longer under the control of a job control
 3968 shell and hence would not normally ever be continued. Because of this, orphaned processes that
 3969 receive terminal-related stop signals (SIGTSTP , SIGTTIN, SIGTTOU, but not SIGSTOP) must
 3970 not be allowed to stop. The goal is to prevent stopped processes from languishing forever. (As
 3971 SIGSTOP is sent only via *kill()*, it is assumed that the process or user sending a SIGSTOP can
 3972 send a SIGCONT when desired.) Instead, the system must discard the stop signal. As an
 3973 extension, it may also deliver another signal in its place. 4.3 BSD sends a SIGKILL, which is
 3974 overly effective because SIGKILL is not catchable. Another possible choice is SIGHUP. 4.3 BSD
 3975 also does this for orphaned processes (processes whose parent has terminated) rather than for
 3976 members of orphaned process groups; this is less desirable because job control shells manage
 3977 process groups. POSIX.1 also prevents SIGTTIN and SIGTTOU signals from being generated for
 3978 processes in orphaned process groups as a direct result of activity on a terminal, preventing
 3979 infinite loops when *read()* and *write()* calls generate signals that are discarded; see Section
 3980 A.11.1.4 (on page 3357). A similar restriction on the generation of SIGTSTP was considered, but
 3981 that would be unnecessary and more difficult to implement due to its asynchronous nature.

3982 Although POSIX.1 requires that signal-catching functions be called with only one argument,
 3983 there is nothing to prevent conforming implementations from extending POSIX.1 to pass
 3984 additional arguments, as long as Strictly Conforming POSIX.1 Applications continue to compile
 3985 and execute correctly. Most historical implementations do, in fact, pass additional, signal-
 3986 specific arguments to certain signal-catching routines.

3987 There was a proposal to change the declared type of the signal handler to:

```
3988     void func (int sig, ...);
```

3989 The usage of ellipses ("...") is ISO C standard syntax to indicate a variable number of
 3990 arguments. Its use was intended to allow the implementation to pass additional information to
 3991 the signal handler in a standard manner.

3992 Unfortunately, this construct would require all signal handlers to be defined with this syntax
 3993 because the ISO C standard allows implementations to use a different parameter passing
 3994 mechanism for variable parameter lists than for non-variable parameter lists. Thus, all existing
 3995 signal handlers in all existing applications would have to be changed to use the variable syntax
 3996 in order to be standard and portable. This is in conflict with the goal of Minimal Changes to
 3997 Existing Application Code.

3998 When terminating a process from a signal-catching function, processes should be aware of any
 3999 interpretation that their parent may make of the status returned by *wait()* or *waitpid()*. In
 4000 particular, a signal-catching function should not call *exit(0)* or *_exit(0)* unless it wants to indicate
 4001 successful termination. A non-zero argument to *exit()* or *_exit()* can be used to indicate
 4002 unsuccessful termination. Alternatively, the process can use *kill()* to send itself a fatal signal
 4003 (first ensuring that the signal is set to the default action and not blocked). See also the
 4004 RATIONALE section of the *_exit()* function.

4005 The behavior of *unsafe* functions, as defined by this section, is undefined when they are invoked
4006 from signal-catching functions in certain circumstances. The behavior of reentrant functions, as
4007 defined by this section, is as specified by POSIX.1, regardless of invocation from a signal-
4008 catching function. This is the only intended meaning of the statement that reentrant functions
4009 may be used in signal-catching functions without restriction. Applications must still consider all
4010 effects of such functions on such things as data structures, files, and process state. In particular,
4011 application writers need to consider the restrictions on interactions when interrupting *sleep()*
4012 (see *sleep()*) and interactions among multiple handles for a file description. The fact that any
4013 specific function is listed as reentrant does not necessarily mean that invocation of that function
4014 from a signal-catching function is recommended.

4015 In order to prevent errors arising from interrupting non-reentrant function calls, applications
4016 should protect calls to these functions either by blocking the appropriate signals or through the
4017 use of some programmatic semaphore. POSIX.1 does not address the more general problem of
4018 synchronizing access to shared data structures. Note in particular that even the “safe” functions
4019 may modify the global variable *errno*; the signal-catching function may want to save and restore
4020 its value. The same principles apply to the reentrancy of application routines and asynchronous
4021 data access.

4022 Note that *longjmp()* and *siglongjmp()* are not in the list of reentrant functions. This is because the
4023 code executing after *longjmp()* or *siglongjmp()* can call any unsafe functions with the same
4024 danger as calling those unsafe functions directly from the signal handler. Applications that use
4025 *longjmp()* or *siglongjmp()* out of signal handlers require rigorous protection in order to be
4026 portable. Many of the other functions that are excluded from the list are traditionally
4027 implemented using either the C language *malloc()* or *free()* functions or the ISO C standard I/O
4028 library, both of which traditionally use data structures in a non-reentrant manner. Because any
4029 combination of different functions using a common data structure can cause reentrancy
4030 problems, POSIX.1 does not define the behavior when any unsafe function is called in a signal
4031 handler that interrupts any unsafe function.

4032 The only realtime extension to signal actions is the addition of the additional parameters to the
4033 signal-catching function. This extension has been explained and motivated in the previous
4034 section. In making this extension, though, developers of POSIX.1b ran into issues relating to
4035 function prototypes. In response to input from the POSIX.1 standard developers, members were
4036 added to the **sigaction** structure to specify function prototypes for the newer signal-catching
4037 function specified by POSIX.1b. These members follow changes that are being made to POSIX.1.
4038 Note that IEEE Std 1003.1-200x explicitly states that these fields may overlap so that a union can
4039 be defined. This will enable existing implementations of POSIX.1 to maintain binary-
4040 compatibility when these extensions are added.

4041 The **siginfo_t** structure was adopted for passing the application-defined value to match existing
4042 practice, but the existing practice has no provision for an application-defined value, so this was
4043 added. Note that POSIX normally reserves the “_t” type designation for opaque types. The
4044 **siginfo_t** structure breaks with this convention to follow existing practice and thus promote
4045 portability. Standardization of the existing practice for the other members of this structure may
4046 be addressed in the future.

4047 Although it is not explicitly visible to applications, there are additional semantics for signal
4048 actions implied by queued signals and their interaction with other POSIX.1b realtime functions.
4049 Specifically:

- 4050 • It is not necessary to queue signals whose action is SIG_IGN.
- 4051 • For implementations that support POSIX.1b timers, some interaction with the timer functions
4052 at signal delivery is implied to manage the timer overrun count.

4053 **B.2.4.4** *Signal Effects on Other Functions*

4054 The most common behavior of an interrupted function after a signal-catching function returns is
4055 for the interrupted function to give an [EINTR] error. However, there are a number of specific
4056 exceptions, including *sleep()* and certain situations with *read()* and *write()*.

4057 The historical implementations of many functions defined by IEEE Std 1003.1-200x are not
4058 interruptible, but delay delivery of signals generated during their execution until after they
4059 complete. This is never a problem for functions that are guaranteed to complete in a short
4060 (imperceptible to a human) period of time. It is normally those functions that can suspend a
4061 process indefinitely or for long periods of time (for example, *wait()*, *pause()*, *sigsuspend()*, *sleep()*,
4062 or *read()/write()* on a slow device like a terminal] that are interruptible. This permits
4063 applications to respond to interactive signals or to set timeouts on calls to most such functions
4064 with *alarm()*. Therefore, implementations should generally make such functions (including ones
4065 defined as extensions) interruptible.

4066 Functions not mentioned explicitly as interruptible may be so on some implementations,
4067 possibly as an extension where the function gives an [EINTR] error. There are several functions
4068 (for example, *getpid()*, *getuid()*) that are specified as never returning an error, which can thus
4069 never be extended in this way.

4070 **B.2.5** **Standard I/O Streams**4071 **B.2.5.1** *Interaction of File Descriptors and Standard I/O Streams*

4072 There is no additional rationale provided for this section.

4073 **B.2.5.2** *Stream Orientation and Encoding Rules*

4074 There is no additional rationale provided for this section.

4075 **B.2.6** **STREAMS**

4076 STREAMS are introduced into IEEE Std 1003.1-200x as part of the alignment with the Single
4077 UNIX Specification, but marked as an option in recognition that not all systems may wish to
4078 implement the facility. The option within IEEE Std 1003.1-200x is denoted by the XSR margin
4079 marker. The standard developers made this option independent of the XSI option.

4080 STREAMS are a method of implementing network services and other character-based
4081 input/output mechanisms, with the STREAM being a full-duplex connection between a process
4082 and a device. STREAMS provides direct access to protocol modules, and optional protocol
4083 modules can be interposed between the process-end of the STREAM and the device-driver at the
4084 device-end of the STREAM. Pipes can be implemented using the STREAMS mechanism, so they
4085 can provide process-to-process as well as process-to-device communications.

4086 This section introduces STREAMS I/O, the message types used to control them, an overview of
4087 the priority mechanism, and the interfaces used to access them.

4088 **B.2.6.1** *Accessing STREAMS*

4089 There is no additional rationale provided for this section.

4090 B.2.7 XSI Interprocess Communication

4091 There are two forms of IPC supported as options in IEEE Std 1003.1-200x. The traditional
4092 System V IPC routines derived from the SVID—that is, the *msg**(), *sem**(), and *shm**()
4093 interfaces—are mandatory on XSI-conformant systems. Thus, all XSI-conformant systems
4094 provide the same mechanisms for manipulating messages, shared memory, and semaphores.

4095 In addition, the POSIX Realtime Extension provides an alternate set of routines for those systems
4096 supporting the appropriate options.

4097 The application writer is presented with a choice: the System V interfaces or the POSIX
4098 interfaces (loosely derived from the Berkeley interfaces). The XSI profile prefers the System V
4099 interfaces, but the POSIX interfaces may be more suitable for realtime or other performance-
4100 sensitive applications.

4101 B.2.7.1 IPC General Information

4102 General information that is shared by all three mechanisms is described in this section. The
4103 common permissions mechanism is briefly introduced, describing the mode bits, and how they
4104 are used to determine whether or not a process has access to read or write/alter the appropriate
4105 instance of one of the IPC mechanisms. All other relevant information is contained in the
4106 reference pages themselves.

4107 The semaphore type of IPC allows processes to communicate through the exchange of
4108 semaphore values. A semaphore is a positive integer. Since many applications require the use of
4109 more than one semaphore, XSI-conformant systems have the ability to create sets or arrays of
4110 semaphores.

4111 Calls to support semaphores include:

4112 *semctl*(), *semget*(), *semop*()

4113 Semaphore sets are created by using the *semget*() function.

4114 The message type of IPC allows process to communicate through the exchange of data stored in
4115 buffers. This data is transmitted between processes in discrete portions known as messages.

4116 Calls to support message queues include:

4117 *msgctl*(), *msgget*(), *msgrcv*(), *msgsnd*()

4118 The share memory type of IPC allows two or more processes to share memory and consequently
4119 the data contained therein. This is done by allowing processes to set up access to a common
4120 memory address space. This sharing of memory provides a fast means of exchange of data
4121 between processes.

4122 Calls to support shared memory include:

4123 *shmctl*(), *shmdt*(), *shmget*()

4124 The *ftok*() interface is also provided.

4125 **B.2.8 Realtime**4126 **Advisory Information**

4127 POSIX.1b contains an Informative Annex with proposed interfaces for “real-time files”. These
 4128 interfaces could determine groups of the exact parameters required to do “direct I/O” or
 4129 “extents”. These interfaces were objected to by a significant portion of the balloting group as too
 4130 complex. A conforming application had little chance of correctly navigating the large parameter
 4131 space to match its desires to the system. In addition, they only applied to a new type of file
 4132 (realtime files) and they told the implementation exactly what to do as opposed to advising the
 4133 implementation on application behavior and letting it optimize for the system the (portable)
 4134 application was running on. For example, it was not clear how a system that had a disk array
 4135 should set its parameters.

4136 There seemed to be several overall goals:

- 4137 • Optimizing sequential access
- 4138 • Optimizing caching behavior
- 4139 • Optimizing I/O data transfer
- 4140 • Preallocation

4141 The advisory interfaces, *posix_fadvise()* and *posix_madvise()*, satisfy the first two goals. The
 4142 `POSIX_FADV_SEQUENTIAL` and `POSIX_MADV_SEQUENTIAL` advice tells the
 4143 implementation to expect serial access. Typically the system will prefetch the next several serial
 4144 accesses in order to overlap I/O. It may also free previously accessed serial data if memory is
 4145 tight. If the application is not doing serial access it can use `POSIX_FADV_WILLNEED` and
 4146 `POSIX_MADV_WILLNEED` to accomplish I/O overlap, as required. When the application
 4147 advises `POSIX_FADV_RANDOM` or `POSIX_MADV_RANDOM` behavior, the implementation
 4148 usually tries to fetch a minimum amount of data with each request and it does not expect much
 4149 locality. `POSIX_FADV_DONTNEED` and `POSIX_MADV_DONTNEED` allow the system to free
 4150 up caching resources as the data will not be required in the near future.

4151 `POSIX_FADV_NOREUSE` tells the system that caching the specified data is not optimal. For file
 4152 I/O, the transfer should go directly to the user buffer instead of being cached internally by the
 4153 implementation. To portably perform direct disk I/O on all systems, the application must
 4154 perform its I/O transfers according to the following rules:

- 4155 1. The user buffer should be aligned according to the `{POSIX_REC_XFER_ALIGN} pathconf()`
 4156 variable.
- 4157 2. The number of bytes transferred in an I/O operation should be a multiple of the
 4158 `{POSIX_ALLOC_SIZE_MIN} pathconf()` variable.
- 4159 3. The offset into the file at the start of an I/O operation should be a multiple of the
 4160 `{POSIX_ALLOC_SIZE_MIN} pathconf()` variable.
- 4161 4. The application should ensure that all threads which open a given file specify
 4162 `POSIX_FADV_NOREUSE` to be sure that there is no unexpected interaction between
 4163 threads using buffered I/O and threads using direct I/O to the same file.

4164 In some cases, a user buffer must be properly aligned in order to be transferred directly to/from
 4165 the device. The `{POSIX_REC_XFER_ALIGN} pathconf()` variable tells the application the proper
 4166 alignment.

4167 The preallocation goal is met by the space control function, *posix_fallocate()*. The application can
 4168 use *posix_fallocate()* to guarantee no [ENOSPC] errors and to improve performance by prepaying

4169 any overhead required for block allocation.

4170 Implementations may use information conveyed by a previous *posix_fadvise()* call to influence
4171 the manner in which allocation is performed. For example, if an application did the following
4172 calls:

```
4173     fd = open("file");
4174     posix_fadvise(fd, offset, len, POSIX_FADV_SEQUENTIAL);
4175     posix_fallocate(fd, len, size);
```

4176 an implementation might allocate the file contiguously on disk.

4177 Finally, the *pathconf()* variables {*POSIX_REC_MIN_XFER_SIZE*},
4178 {*POSIX_REC_MAX_XFER_SIZE*}, and {*POSIX_REC_INCR_XFER_SIZE*} tell the application a
4179 range of transfer sizes that are recommended for best I/O performance.

4180 Where bounded response time is required, the vendor can supply the appropriate settings of the
4181 advisories to achieve a guaranteed performance level.

4182 The interfaces meet the goals while allowing applications using regular files to take advantage of
4183 performance optimizations. The interfaces tell the implementation expected application
4184 behavior which the implementation can use to optimize performance on a particular system
4185 with a particular dynamic load.

4186 The *posix_memalign()* function was added to allow for the allocation of specifically aligned
4187 buffers; for example, for {*POSIX_REC_XFER_ALIGN*}.

4188 The working group also considered the alternative of adding a function which would return an
4189 aligned pointer to memory within a user supplied buffer. This was not considered to be the best
4190 method, because it potentially wastes large amounts of memory when buffers need to be aligned
4191 on large alignment boundaries.

4192 **Message Passing**

4193 This section provides the rationale for the definition of the message passing interface in
4194 IEEE Std 1003.1-200x. This is presented in terms of the objectives, models, and requirements
4195 imposed upon this interface.

4196 • Objectives

4197 Many applications, including both realtime and database applications, require a means of
4198 passing arbitrary amounts of data between cooperating processes comprising the overall
4199 application on one or more processors. Many conventional interfaces for interprocess
4200 communication are insufficient for realtime applications in that efficient and deterministic
4201 data passing methods cannot be implemented. This has prompted the definition of message
4202 passing interfaces providing these facilities:

- 4203 — Open a message queue.
- 4204 — Send a message to a message queue.
- 4205 — Receive a message from a queue, either synchronously or asynchronously.
- 4206 — Alter message queue attributes for flow and resource control.

4207 It is assumed that an application may consist of multiple cooperating processes and that
4208 these processes may wish to communicate and coordinate their activities. The message
4209 passing facility described in IEEE Std 1003.1-200x allows processes to communicate through
4210 system-wide queues. These message queues are accessed through names that may be
4211 pathnames. A message queue can be opened for use by multiple sending and/or multiple

- 4212 receiving processes.
- 4213 • Background on Embedded Applications
- 4214 Interprocess communication utilizing message passing is a key facility for the construction of
4215 deterministic, high-performance realtime applications. The facility is present in all realtime
4216 systems and is the framework upon which the application is constructed. The performance of
4217 the facility is usually a direct indication of the performance of the resulting application.
- 4218 Realtime applications, especially for embedded systems, are typically designed around the
4219 performance constraints imposed by the message passing mechanisms. Applications for
4220 embedded systems are typically very tightly constrained. Application writers expect to
4221 design and control the entire system. In order to minimize system costs, the writer will
4222 attempt to use all resources to their utmost and minimize the requirement to add additional
4223 memory or processors.
- 4224 The embedded applications usually share address spaces and only a simple message passing
4225 mechanism is required. The application can readily access common data incurring only
4226 mutual-exclusion overheads. The models desired are the simplest possible with the
4227 application building higher-level facilities only when needed.
- 4228 • Requirements
- 4229 The following requirements determined the features of the message passing facilities defined
4230 in IEEE Std 1003.1-200x:
- 4231 — Naming of Message Queues
- 4232 The mechanism for gaining access to a message queue is a pathname evaluated in a
4233 context that is allowed to be a file system name space, or it can be independent of any file
4234 system. This is a specific attempt to allow implementations based on either method in
4235 order to address both embedded systems and to also allow implementation in larger
4236 systems.
- 4237 The interface of *mq_open()* is defined to allow but not require the access control and name
4238 conflicts resulting from utilizing a file system for name resolution. All required behavior
4239 is specified for the access control case. Yet a conforming implementation, such as an
4240 embedded system kernel, may define that there are no distinctions between users and
4241 may define that all process have all access privileges.
- 4242 — Embedded System Naming
- 4243 Embedded systems need to be able to utilize independent name spaces for accessing the
4244 various system objects. They typically do not have a file system, precluding its utilization
4245 as a common name resolution mechanism. The modularity of an embedded system limits
4246 the connections between separate mechanisms that can be allowed.
- 4247 Embedded systems typically do not have any access protection. Since the system does not
4248 support the mixing of applications from different areas, and usually does not even have
4249 the concept of an authorization entity, access control is not useful.
- 4250 — Large System Naming
- 4251 On systems with more functionality, the name resolution must support the ability to use
4252 the file system as the name resolution mechanism/object storage medium and to have
4253 control over access to the objects. Utilizing the pathname space can result in further errors
4254 when the names conflict with other objects.
- 4255 — Fixed Size of Messages

4256 The interfaces impose a fixed upper bound on the size of messages that can be sent to a
 4257 specific message queue. The size is set on an individual queue basis and cannot be
 4258 changed dynamically.

4259 The purpose of the fixed size is to increase the ability of the system to optimize the
 4260 implementation of *mq_send()* and *mq_receive()*. With fixed sizes of messages and fixed
 4261 numbers of messages, specific message blocks can be pre-allocated. This eliminates a
 4262 significant amount of checking for errors and boundary conditions. Additionally, an
 4263 implementation can optimize data copying to maximize performance. Finally, with a
 4264 restricted range of message sizes, an implementation is better able to provide
 4265 deterministic operations.

4266 — Prioritization of Messages

4267 Message prioritization allows the application to determine the order in which messages
 4268 are received. Prioritization of messages is a key facility that is provided by most realtime
 4269 kernels and is heavily utilized by the applications. The major purpose of having priorities
 4270 in message queues is to avoid priority inversions in the message system, where a high-
 4271 priority message is delayed behind one or more lower-priority messages. This allows the
 4272 applications to be designed so that they do not need to be interrupted in order to change
 4273 the flow of control when exceptional conditions occur. The prioritization does add
 4274 additional overhead to the message operations in those cases it is actually used but a
 4275 clever implementation can optimize for the FIFO case to make that more efficient.

4276 — Asynchronous Notification

4277 The interface supports the ability to have a task asynchronously notified of the
 4278 availability of a message on the queue. The purpose of this facility is to allow the task to
 4279 perform other functions and yet still be notified that a message has become available on
 4280 the queue.

4281 To understand the requirement for this function, it is useful to understand two models of
 4282 application design: a single task performing multiple functions and multiple tasks
 4283 performing a single function. Each of these models has advantages.

4284 Asynchronous notification is required to build the model of a single task performing
 4285 multiple operations. This model typically results from either the expectation that
 4286 interruption is less expensive than utilizing a separate task or from the growth of the
 4287 application to include additional functions.

4288 **Semaphores**

4289 Semaphores are a high-performance process synchronization mechanism. Semaphores are
 4290 named by null-terminated strings of characters.

4291 A semaphore is created using the *sem_init()* function or the *sem_open()* function with the
 4292 *O_CREAT* flag set in *oflag*.

4293 To use a semaphore, a process has to first initialize the semaphore or inherit an open descriptor
 4294 for the semaphore via *fork()*.

4295 A semaphore preserves its state when the last reference is closed. For example, if a semaphore
 4296 has a value of 13 when the last reference is closed, it will have a value of 13 when it is next
 4297 opened.

4298 When a semaphore is created, an initial state for the semaphore has to be provided. This value is
 4299 a non-negative integer. Negative values are not possible since they indicate the presence of
 4300 blocked processes. The persistence of any of these objects across a system crash or a system

- 4301 reboot is undefined. Conforming applications shall not depend on any sort of persistence across
4302 a system reboot or a system crash.
- 4303 • Models and Requirements
- 4304 A realtime system requires synchronization and communication between the processes
4305 comprising the overall application. An efficient and reliable synchronization mechanism has
4306 to be provided in a realtime system that will allow more than one schedulable process
4307 mutually-exclusive access to the same resource. This synchronization mechanism has to
4308 allow for the optimal implementation of synchronization or systems implementors will
4309 define other, more cost-effective methods.
- 4310 At issue are the methods whereby multiple processes (tasks) can be designed and
4311 implemented to work together in order to perform a single function. This requires
4312 interprocess communication and synchronization. A semaphore mechanism is the lowest
4313 level of synchronization that can be provided by an operating system.
- 4314 A semaphore is defined as an object that has an integral value and a set of blocked processes
4315 associated with it. If the value is positive or zero, then the set of blocked processes is empty;
4316 otherwise, the size of the set is equal to the absolute value of the semaphore value. The value
4317 of the semaphore can be incremented or decremented by any process with access to the
4318 semaphore and must be done as an indivisible operation. When a semaphore value is less
4319 than or equal to zero, any process that attempts to lock it again will block or be informed that
4320 it is not possible to perform the operation.
- 4321 A semaphore may be used to guard access to any resource accessible by more than one
4322 schedulable task in the system. It is a global entity and not associated with any particular
4323 process. As such, a method of obtaining access to the semaphore has to be provided by the
4324 operating system. A process that wants access to a critical resource (section) has to wait on
4325 the semaphore that guards that resource. When the semaphore is locked on behalf of a
4326 process, it knows that it can utilize the resource without interference by any other
4327 cooperating process in the system. When the process finishes its operation on the resource,
4328 leaving it in a well-defined state, it posts the semaphore, indicating that some other process
4329 may now obtain the resource associated with that semaphore.
- 4330 In this section, mutexes and condition variables are specified as the synchronization
4331 mechanisms between threads.
- 4332 These primitives are typically used for synchronizing threads that share memory in a single
4333 process. However, this section provides an option allowing the use of these synchronization
4334 interfaces and objects between processes that share memory, regardless of the method for
4335 sharing memory.
- 4336 Much experience with semaphores shows that there are two distinct uses of synchronization:
4337 locking, which is typically of short duration; and waiting, which is typically of long or
4338 unbounded duration. These distinct usages map directly onto mutexes and condition
4339 variables, respectively.
- 4340 Semaphores are provided in IEEE Std 1003.1-200x primarily to provide a means of
4341 synchronization for processes; these processes may or may not share memory. Mutexes and
4342 condition variables are specified as synchronization mechanisms between threads; these
4343 threads always share (some) memory. Both are synchronization paradigms that have been in
4344 widespread use for a number of years. Each set of primitives is particularly well matched to
4345 certain problems.
- 4346 With respect to binary semaphores, experience has shown that condition variables and
4347 mutexes are easier to use for many synchronization problems than binary semaphores. The

4348 primary reason for this is the explicit appearance of a Boolean predicate that specifies when
4349 the condition wait is satisfied. This Boolean predicate terminates a loop, including the call to
4350 *pthread_cond_wait()*. As a result, extra wakeups are benign since the predicate governs
4351 whether the thread will actually proceed past the condition wait. With stateful primitives,
4352 such as binary semaphores, the wakeup in itself typically means that the wait is satisfied. The
4353 burden of ensuring correctness for such waits is thus placed on *all* signalers of the semaphore
4354 rather than on an *explicitly coded* Boolean predicate located at the condition wait. Experience
4355 has shown that the latter creates a major improvement in safety and ease-of-use.

4356 Counting semaphores are well matched to dealing with producer/consumer problems,
4357 including those that might exist between threads of different processes, or between a signal
4358 handler and a thread. In the former case, there may be little or no memory shared by the
4359 processes; in the latter case, one is not communicating between co-equal threads, but
4360 between a thread and an interruptlike entity. It is for these reasons that IEEE Std 1003.1-200x
4361 allows semaphores to be used by threads.

4362 Mutexes and condition variables have been effectively used with and without priority
4363 inheritance, priority ceiling, and other attributes to synchronize threads that share memory.
4364 The efficiency of their implementation is comparable to or better than that of other
4365 synchronization primitives that are sometimes harder to use (for example, binary
4366 semaphores). Furthermore, there is at least one known implementation of Ada tasking that
4367 uses these primitives. Mutexes and condition variables together constitute an appropriate,
4368 sufficient, and complete set of interthread synchronization primitives.

4369 Efficient multi-threaded applications require high-performance synchronization primitives.
4370 Considerations of efficiency and generality require a small set of primitives upon which more
4371 sophisticated synchronization functions can be built.

4372 • Standardization Issues

4373 It is possible to implement very high-performance semaphores using test-and-set
4374 instructions on shared memory locations. The library routines that implement such a high-
4375 performance interface has to properly ensure that a *sem_wait()* or *sem_trywait()* operation
4376 that cannot be performed will issue a blocking semaphore system call or properly report the
4377 condition to the application. The same interface to the application program would be
4378 provided by a high-performance implementation.

4379 *B.2.8.1 Realtime Signals*

4380 **Realtime Signals Extension**

4381 This portion of the rationale presents models, requirements, and standardization issues relevant
4382 to the Realtime Signals Extension. This extension provides the capability required to support
4383 reliable, deterministic, asynchronous notification of events. While a new mechanism,
4384 unencumbered by the historical usage and semantics of POSIX.1 signals, might allow for a more
4385 efficient implementation, the application requirements for event notification can be met with a
4386 small number of extensions to signals. Therefore, a minimal set of extensions to signals to
4387 support the application requirements is specified.

4388 The realtime signal extensions specified in this section are used by other realtime functions
4389 requiring asynchronous notification:

4390 • Models

4391 The model supported is one of multiple cooperating processes, each of which handles
4392 multiple asynchronous external events. Events represent occurrences that are generated as

- 4393 the result of some activity in the system. Examples of occurrences that can constitute an
4394 event include:
- 4395 — Completion of an asynchronous I/O request
 - 4396 — Expiration of a POSIX.1b timer
 - 4397 — Arrival of an interprocess message
 - 4398 — Generation of a user-defined event
- 4399 Processing of these events may occur synchronously via polling for event notifications or
4400 asynchronously via a software interrupt mechanism. Existing practice for this model is well
4401 established for traditional proprietary realtime operating systems, realtime executives, and
4402 realtime extended POSIX-like systems.
- 4403 A contrasting model is that of “cooperating sequential processes” where each process
4404 handles a single priority of events via polling. Each process blocks while waiting for events,
4405 and each process depends on the preemptive, priority-based process scheduling mechanism
4406 to arbitrate between events of different priority that need to be processed concurrently.
4407 Existing practice for this model is also well established for small realtime executives that
4408 typically execute in an unprotected physical address space, but it is just emerging in the
4409 context of a fuller function operating system with multiple virtual address spaces.
- 4410 It could be argued that the cooperating sequential process model, and the facilities supported
4411 by the POSIX Threads Extension obviate a software interrupt model. But, even with the
4412 cooperating sequential process model, the need has been recognized for a software interrupt
4413 model to handle exceptional conditions and process aborting, so the mechanism must be
4414 supported in any case. Furthermore, it is not the purview of IEEE Std 1003.1-200x to attempt
4415 to convince realtime practitioners that their current application models based on software
4416 interrupts are “broken” and should be replaced by the cooperating sequential process model.
4417 Rather, it is the charter of IEEE Std 1003.1-200x to provide standard extensions to
4418 mechanisms that support existing realtime practice.
- 4419 • Requirements
- 4420 This section discusses the following realtime application requirements for asynchronous
4421 event notification:
- 4422 — Reliable delivery of asynchronous event notification
- 4423 The events notification mechanism shall guarantee delivery of an event notification.
4424 Asynchronous operations (such as asynchronous I/O and timers) that complete
4425 significantly after they are invoked have to guarantee that delivery of the event
4426 notification can occur at the time of completion.
- 4427 — Prioritized handling of asynchronous event notifications
- 4428 The events notification mechanism shall support the assigning of a user function as an
4429 event notification handler. Furthermore, the mechanism shall support the preemption of
4430 an event handler function by a higher priority event notification and shall support the
4431 selection of the highest priority pending event notification when multiple notifications (of
4432 different priority) are pending simultaneously.
- 4433 The model here is based on hardware interrupts. Asynchronous event handling allows
4434 the application to ensure that time-critical events are immediately processed when
4435 delivered, without the indeterminism of being at a random location within a polling loop.
4436 Use of handler priority allows the specification of how handlers are interrupted by other
4437 higher priority handlers.

- 4438 — Differentiation between multiple occurrences of event notifications of the same type
- 4439 The events notification mechanism shall pass an application-defined value to the event
4440 handler function. This value can be used for a variety of purposes, such as enabling the
4441 application to identify which of several possible events of the same type (for example,
4442 timer expirations) has occurred.
- 4443 — Polled reception of asynchronous event notifications
- 4444 The events notification mechanism shall support blocking and non-blocking polls for
4445 asynchronous event notification.
- 4446 The polled mode of operation is often preferred over the interrupt mode by those
4447 practitioners accustomed to this model. Providing support for this model facilitates the
4448 porting of applications based on this model to POSIX.1b conforming systems.
- 4449 — Deterministic response to asynchronous event notifications
- 4450 The events notification mechanism shall not preclude implementations that provide
4451 deterministic event dispatch latency and shall minimize the number of system calls
4452 needed to use the event facilities during realtime processing.
- 4453 • Rationale for Extension
- 4454 POSIX.1 signals have many of the characteristics necessary to support the asynchronous
4455 handling of event notifications, and the Realtime Signals Extension addresses the following
4456 deficiencies in the POSIX.1 signal mechanism:
- 4457 — Signals do not support reliable delivery of event notification. Subsequent occurrences of
4458 a pending signal are not guaranteed to be delivered.
- 4459 — Signals do not support prioritized delivery of event notifications. The order of signal
4460 delivery when multiple unblocked signals are pending is undefined.
- 4461 — Signals do not support the differentiation between multiple signals of the same type.

4462 *B.2.8.2 Asynchronous I/O*

4463 Many applications need to interact with the I/O subsystem in an asynchronous manner. The
4464 asynchronous I/O mechanism provides the ability to overlap application processing and I/O
4465 operations initiated by the application. The asynchronous I/O mechanism allows a single
4466 process to perform I/O simultaneously to a single file multiple times or to multiple files
4467 multiple times.

4468 **Overview**

4469 Asynchronous I/O operations proceed in logical parallel with the processing done by the
4470 application after the asynchronous I/O has been initiated. Other than this difference,
4471 asynchronous I/O behaves similarly to normal I/O using *read()*, *write()*, *lseek()*, and *fsync()*.
4472 The effect of issuing an asynchronous I/O request is as if a separate thread of execution were to
4473 perform atomically the implied *lseek()* operation, if any, and then the requested I/O operation
4474 (either *read()*, *write()*, or *fsync()*). There is no seek implied with a call to *aio_fsync()*. Concurrent
4475 asynchronous operations and synchronous operations applied to the same file update the file as
4476 if the I/O operations had proceeded serially.

4477 When asynchronous I/O completes, a signal can be delivered to the application to indicate the
4478 completion of the I/O. This signal can be used to indicate that buffers and control blocks used
4479 for asynchronous I/O can be reused. Signal delivery is not required for an asynchronous
4480 operation and may be turned off on a per-operation basis by the application. Signals may also be

4481 synchronously polled using *aio_suspend()*, *sigtimedwait()*, or *sigwaitinfo()*.

4482 Normal I/O has a return value and an error status associated with it. Asynchronous I/O returns
4483 a value and an error status when the operation is first submitted, but that only relates to whether
4484 the operation was successfully queued up for servicing. The I/O operation itself also has a
4485 return status and an error value. To allow the application to retrieve the return status and the
4486 error value, functions are provided that, given the address of an asynchronous I/O control
4487 block, yield the return and error status associated with the operation. Until an asynchronous I/O
4488 operation is done, its error status shall be [EINPROGRESS]. Thus, an application can poll for
4489 completion of an asynchronous I/O operation by waiting for the error status to become equal to
4490 a value other than [EINPROGRESS]. The return status of an asynchronous I/O operation is
4491 undefined so long as the error status is equal to [EINPROGRESS].

4492 Storage for asynchronous operation return and error status may be limited. Submission of
4493 asynchronous I/O operations may fail if this storage is exceeded. When an application retrieves
4494 the return status of a given asynchronous operation, therefore, any system-maintained storage
4495 used for this status and the error status may be reclaimed for use by other asynchronous
4496 operations.

4497 Asynchronous I/O can be performed on file descriptors that have been enabled for POSIX.1b
4498 synchronized I/O. In this case, the I/O operation still occurs asynchronously, as defined herein;
4499 however, the asynchronous operation I/O in this case is not completed until the I/O has reached
4500 either the state of synchronized I/O data integrity completion or synchronized I/O file integrity
4501 completion, depending on the sort of synchronized I/O that is enabled on the file descriptor.

4502 **Models**

4503 Three models illustrate the use of asynchronous I/O: a journalization model, a data acquisition
4504 model, and a model of the use of asynchronous I/O in supercomputing applications.

- 4505 • **Journalization Model**

4506 Many realtime applications perform low-priority journalizing functions. Journalizing
4507 requires that logging records be queued for output without blocking the initiating process.

- 4508 • **Data Acquisition Model**

4509 A data acquisition process may also serve as a model. The process has two or more channels
4510 delivering intermittent data that must be read within a certain time. The process issues one
4511 asynchronous read on each channel. When one of the channels needs data collection, the
4512 process reads the data and posts it through an asynchronous write to secondary memory for
4513 future processing.

- 4514 • **Supercomputing Model**

4515 The supercomputing community has used asynchronous I/O much like that specified herein
4516 for many years. This community requires the ability to perform multiple I/O operations to
4517 multiple devices with a minimal number of entries to “the system”; each entry to “the
4518 system” provokes a major delay in operations when compared to the normal progress made
4519 by the application. This existing practice motivated the use of combined *lseek()* and *read()* or
4520 *write()* calls, as well as the *lio_listio()* call. Another common practice is to disable signal
4521 notification for I/O completion, and simply poll for I/O completion at some interval by
4522 which the I/O should be completed. Likewise, interfaces like *aio_cancel()* have been in
4523 successful commercial use for many years. Note also that an underlying implementation of
4524 asynchronous I/O will require the ability, at least internally, to cancel outstanding
4525 asynchronous I/O, at least when the process exits. (Consider an asynchronous read from a
4526 terminal, when the process intends to exit immediately.)

4527 **Requirements**

4528 Asynchronous input and output for realtime implementations have these requirements:

- 4529 • The ability to queue multiple asynchronous read and write operations to a single open
4530 instance. Both sequential and random access should be supported.
- 4531 • The ability to queue asynchronous read and write operations to multiple open instances.
- 4532 • The ability to obtain completion status information by polling and/or asynchronous event
4533 notification.
- 4534 • Asynchronous event notification on asynchronous I/O completion is optional.
- 4535 • It has to be possible for the application to associate the event with the *aiocbp* for the operation
4536 that generated the event.
- 4537 • The ability to cancel queued requests.
- 4538 • The ability to wait upon asynchronous I/O completion in conjunction with other types of
4539 events.
- 4540 • The ability to accept an *aio_read()* and an *aio_cancel()* for a device that accepts a *read()*, and
4541 the ability to accept an *aio_write()* and an *aio_cancel()* for a device that accepts a *write()*. This
4542 does not imply that the operation is asynchronous.

4543 **Standardization Issues**

4544 The following issues are addressed by the standardization of asynchronous I/O:

4545 • Rationale for New Interface

4546 Non-blocking I/O does not satisfy the needs of either realtime or high-performance
4547 computing models; these models require that a process overlap program execution and I/O
4548 processing. Realtime applications will often make use of direct I/O to or from the address
4549 space of the process, or require synchronized (unbuffered) I/O; they also require the ability
4550 to overlap this I/O with other computation. In addition, asynchronous I/O allows an
4551 application to keep a device busy at all times, possibly achieving greater throughput.
4552 Supercomputing and database architectures will often have specialized hardware that can
4553 provide true asynchrony underlying the logical asynchrony provided by this interface. In
4554 addition, asynchronous I/O should be supported by all types of files and devices in the same
4555 manner.

4556 • Effect of Buffering

4557 If asynchronous I/O is performed on a file that is buffered prior to being actually written to
4558 the device, it is possible that asynchronous I/O will offer no performance advantage over
4559 normal I/O; the cycles *stolen* to perform the asynchronous I/O will be taken away from the
4560 running process and the I/O will occur at interrupt time. This potential lack of gain in
4561 performance in no way obviates the need for asynchronous I/O by realtime applications,
4562 which very often will use specialized hardware support; multiple processors; and/or
4563 unbuffered, synchronized I/O.

4564 **B.2.8.3** *Memory Management*

4565 All memory management and shared memory definitions are located in the `<sys/mman.h>`
4566 header. This is for alignment with historical practice.

4567 **Memory Locking Functions**

4568 This portion of the rationale presents models, requirements, and standardization issues relevant
4569 to process memory locking.

4570 • Models

4571 Realtime systems that conform to IEEE Std 1003.1-200x are expected (and desired) to be
4572 supported on systems with demand-paged virtual memory management, non-paged
4573 swapping memory management, and physical memory systems with no memory
4574 management hardware. The general case, however, is the demand-paged, virtual memory
4575 system with each POSIX process running in a virtual address space. Note that this includes
4576 architectures where each process resides in its own virtual address space and architectures
4577 where the address space of each process is only a portion of a larger global virtual address
4578 space.

4579 The concept of memory locking is introduced to eliminate the indeterminacy introduced by
4580 paging and swapping, and to support an upper bound on the time required to access the
4581 memory mapped into the address space of a process. Ideally, this upper bound will be the
4582 same as the time required for the processor to access “main memory”, including any address
4583 translation and cache miss overheads. But some implementations—primarily on
4584 mainframes—will not actually force locked pages to be loaded and held resident in main
4585 memory. Rather, they will handle locked pages so that accesses to these pages will meet the
4586 performance metrics for locked process memory in the implementation. Also, although it is
4587 not, for example, the intention that this interface, as specified, be used to lock process
4588 memory into “cache”, it is conceivable that an implementation could support a large static
4589 RAM memory and define this as “main memory” and use a large[r] dynamic RAM as
4590 “backing store”. These interfaces could then be interpreted as supporting the locking of
4591 process memory into the static RAM. Support for multiple levels of backing store would
4592 require extensions to these interfaces.

4593 Implementations may also use memory locking to guarantee a fixed translation between
4594 virtual and physical addresses where such is beneficial to improving determinacy for
4595 direct-to/from-process input/output. IEEE Std 1003.1-200x does not guarantee to the
4596 application that the virtual-to-physical address translations, if such exist, are fixed, because
4597 such behavior would not be implementable on all architectures on which implementations of
4598 IEEE Std 1003.1-200x are expected. But IEEE Std 1003.1-200x does mandate that an
4599 implementation define, for the benefit of potential users, whether or not locking guarantees
4600 fixed translations.

4601 Memory locking is defined with respect to the address space of a process. Only the pages
4602 mapped into the address space of a process may be locked by the process, and when the
4603 pages are no longer mapped into the address space—for whatever reason—the locks
4604 established with respect to that address space are removed. Shared memory areas warrant
4605 special mention, as they may be mapped into more than one address space or mapped more
4606 than once into the address space of a process; locks may be established on pages within these
4607 areas with respect to several of these mappings. In such a case, the lock state of the
4608 underlying physical pages is the logical OR of the lock state with respect to each of the
4609 mappings. Only when all such locks have been removed are the shared pages considered
4610 unlocked.

4611 In recognition of the page granularity of Memory Management Units (MMU), and in order to
4612 support locking of ranges of address space, memory locking is defined in terms of “page”
4613 granularity. That is, for the interfaces that support an address and size specification for the
4614 region to be locked, the address must be on a page boundary, and all pages mapped by the
4615 specified range are locked, if valid. This means that the length is implicitly rounded up to a
4616 multiple of the page size. The page size is implementation-defined and is available to
4617 applications as a compile time symbolic constant or at runtime via *sysconf()*.

4618 A “real memory” POSIX.1b implementation that has no MMU could elect not to support
4619 these interfaces, returning [ENOSYS]. But an application could easily interpret this as
4620 meaning that the implementation would unconditionally page or swap the application when
4621 such is not the case. It is the intention of IEEE Std 1003.1-200x that such a system could define
4622 these interfaces as “NO-OPs”, returning success without actually performing any function
4623 except for mandated argument checking.

4624 • Requirements

4625 For realtime applications, memory locking is generally considered to be required as part of
4626 application initialization. This locking is performed after an application has been loaded (that
4627 is, *exec'd*) and the program remains locked for its entire lifetime. But to support applications
4628 that undergo major mode changes where, in one mode, locking is required, but in another it
4629 is not, the specified interfaces allow repeated locking and unlocking of memory within the
4630 lifetime of a process.

4631 When a realtime application locks its address space, it should not be necessary for the
4632 application to then “touch” all of the pages in the address space to guarantee that they are
4633 resident or else suffer potential paging delays the first time the page is referenced. Thus,
4634 IEEE Std 1003.1-200x requires that the pages locked by the specified interfaces be resident
4635 when the locking functions return successfully.

4636 Many architectures support system-managed stacks that grow automatically when the
4637 current extent of the stack is exceeded. A realtime application has a requirement to be able to
4638 “preallocate” sufficient stack space and lock it down so that it will not suffer page faults to
4639 grow the stack during critical realtime operation. There was no consensus on a portable way
4640 to specify how much stack space is needed, so IEEE Std 1003.1-200x supports no specific
4641 interface for preallocating stack space. But an application can portably lock down a specific
4642 amount of stack space by specifying *MCL_FUTURE* in a call to *memlockall()* and then calling
4643 a dummy function that declares an automatic array of the desired size.

4644 Memory locking for realtime applications is also generally considered to be an “all or
4645 nothing” proposition. That is, the entire process, or none, is locked down. But, for
4646 applications that have well-defined sections that need to be locked and others that do not,
4647 IEEE Std 1003.1-200x supports an optional set of interfaces to lock or unlock a range of
4648 process addresses. Reasons for locking down a specific range include:

4649 — An asynchronous event handler function that must respond to external events in a
4650 deterministic manner such that page faults cannot be tolerated

4651 — An input/output “buffer” area that is the target for direct-to-process I/O, and the
4652 overhead of implicit locking and unlocking for each I/O call cannot be tolerated

4653 Finally, locking is generally viewed as an “application-wide” function. That is, the
4654 application is globally aware of which regions are locked and which are not over time. This is
4655 in contrast to a function that is used temporarily within a “third party” library routine whose
4656 function is unknown to the application, and therefore must have no “side effects”. The
4657 specified interfaces, therefore, do not support “lock stacking” or “lock nesting” within a
4658 process. But, for pages that are shared between processes or mapped more than once into a

4659 process address space, “lock stacking” is essentially mandated by the requirement that
4660 unlocking of pages that are mapped by more than one process or more than once by the same
4661 process does not affect locks established on the other mappings.

4662 There was some support for “lock stacking” so that locking could be transparently used in
4663 functions or opaque modules. But the consensus was not to burden all implementations with
4664 lock stacking (and reference counting), and an implementation option was proposed. There
4665 were strong objections to the option because applications would have to support both
4666 options in order to remain portable. The consensus was to eliminate lock stacking altogether,
4667 primarily through overwhelming support for the System V “m[un]lock[all]” interface on
4668 which IEEE Std 1003.1-200x is now based.

4669 Locks are not inherited across *fork()*s because some implementations implement *fork()* by
4670 creating new address spaces for the child. In such an implementation, requiring locks to be
4671 inherited would lead to new situations in which a fork would fail due to the inability of the
4672 system to lock sufficient memory to lock both the parent and the child. The consensus was
4673 that there was no benefit to such inheritance. Note that this does not mean that locks are
4674 removed when, for instance, a thread is created in the same address space.

4675 Similarly, locks are not inherited across *exec* because some implementations implement *exec*
4676 by unmapping all of the pages in the address space (which, by definition, removes the locks
4677 on these pages), and maps in pages of the *exec*'d image. In such an implementation, requiring
4678 locks to be inherited would lead to new situations in which *exec* would fail. Reporting this
4679 failure would be very cumbersome to detect in time to report to the calling process, and no
4680 appropriate mechanism exists for informing the *exec*'d process of its status.

4681 It was determined that, if the newly loaded application required locking, it was the
4682 responsibility of that application to establish the locks. This is also in keeping with the
4683 general view that it is the responsibility of the application to be aware of all locks that are
4684 established.

4685 There was one request to allow (not mandate) locks to be inherited across *fork()*, and a
4686 request for a flag, MCL_INHERIT, that would specify inheritance of memory locks across
4687 *exec*s. Given the difficulties raised by this and the general lack of support for the feature in
4688 IEEE Std 1003.1-200x, it was not added. IEEE Std 1003.1-200x does not preclude an
4689 implementation from providing this feature for administrative purposes, such as a “run”
4690 command that will lock down and execute specified program. Additionally, the rationale for
4691 the objection equated *fork()* with creating a thread in the address space. IEEE Std 1003.1-200x
4692 does not mandate releasing locks when creating additional threads in an existing process.

4693 • Standardization Issues

4694 One goal of IEEE Std 1003.1-200x is to define a set of primitives that provide the necessary
4695 functionality for realtime applications, with consideration for the needs of other application
4696 domains where such were identified, which is based to the extent possible on existing
4697 industry practice.

4698 The Memory Locking option is required by many realtime applications to tune performance.
4699 Such a facility is accomplished by placing constraints on the virtual memory system to limit
4700 paging of time of the process or of critical sections of the process. This facility should not be
4701 used by most non-realtime applications.

4702 Optional features provided in IEEE Std 1003.1-200x allow applications to lock selected
4703 address ranges with the caveat that the process is responsible for being aware of the page
4704 granularity of locking and the un-nested nature of the locks.

4705 **Mapped Files Functions**

4706 The Memory Mapped Files option provides a mechanism that allows a process to access files by
4707 directly incorporating file data into its address space. Once a file is “mapped” into a process
4708 address space, the data can be manipulated by instructions as memory. The use of mapped files
4709 can significantly reduce I/O data movement since file data does not have to be copied into
4710 process data buffers as in *read()* and *write()*. If more than one process maps a file, its contents
4711 are shared among them. This provides a low overhead mechanism by which processes can
4712 synchronize and communicate.

4713 • Historical Perspective

4714 Realtime applications have historically been implemented using a collection of cooperating
4715 processes or tasks. In early systems, these processes ran on bare hardware (that is, without an
4716 operating system) with no memory relocation or protection. The application paradigms that
4717 arose from this environment involve the sharing of data between the processes.

4718 When realtime systems were implemented on top of vendor-supplied operating systems, the
4719 paradigm or performance benefits of direct access to data by multiple processes was still
4720 deemed necessary. As a result, operating systems that claim to support realtime applications
4721 must support the shared memory paradigm.

4722 Additionally, a number of realtime systems provide the ability to map specific sections of the
4723 physical address space into the address space of a process. This ability is required if an
4724 application is to obtain direct access to memory locations that have specific properties (for
4725 example, refresh buffers or display devices, dual ported memory locations, DMA target
4726 locations). The use of this ability is common enough to warrant some degree of
4727 standardization of its interface. This ability overlaps the general paradigm of shared
4728 memory in that, in both instances, common global objects are made addressable by
4729 individual processes or tasks.

4730 Finally, a number of systems also provide the ability to map process addresses to files. This
4731 provides both a general means of sharing persistent objects, and using files in a manner that
4732 optimizes memory and swapping space usage.

4733 Simple shared memory is clearly a special case of the more general file mapping capability.
4734 In addition, there is relatively widespread agreement and implementation of the file
4735 mapping interface. In these systems, many different types of objects can be mapped (for
4736 example, files, memory, devices, and so on) using the same mapping interfaces. This
4737 approach both minimizes interface proliferation and maximizes the generality of programs
4738 using the mapping interfaces.

4739 • Memory Mapped Files Usage

4740 A memory object can be concurrently mapped into the address space of one or more
4741 processes. The *mmap()* and *munmap()* functions allow a process to manipulate their address
4742 space by mapping portions of memory objects into it and removing them from it. When
4743 multiple processes map the same memory object, they can share access to the underlying
4744 data. Implementations may restrict the size and alignment of mappings to be on *page*-size
4745 boundaries. The page size, in bytes, is the value of the system-configurable variable
4746 {PAGESIZE}, typically accessed by calling *sysconf()* with a *name* argument of
4747 _SC_PAGESIZE. If an implementation has no restrictions on size or alignment, it may
4748 specify a 1-byte page size.

4749 To map memory, a process first opens a memory object. The *ftruncate()* function can be used
4750 to contract or extend the size of the memory object even when the object is currently
4751 mapped. If the memory object is extended, the contents of the extended areas are zeros.

4752 After opening a memory object, the application maps the object into its address space using
4753 the *mmap()* function call. Once a mapping has been established, it remains mapped until
4754 unmapped with *munmap()*, even if the memory object is closed. The *mprotect()* function can
4755 be used to change the memory protections initially established by *mmap()*.

4756 A *close()* of the file descriptor, while invalidating the file descriptor itself, does not unmap
4757 any mappings established for the memory object. The address space, including all mapped
4758 regions, is inherited on *fork()*. The entire address space is unmapped on process termination
4759 or by successful calls to any of the *exec* family of functions.

4760 The *msync()* function is used to force mapped file data to permanent storage.

4761 • Effects on Other Functions

4762 When the Memory Mapped Files option is supported, the operation of the *open()*, *creat()*, and
4763 *unlink()* functions are a natural result of using the file system name space to map the global
4764 names for memory objects.

4765 The *ftruncate()* function can be used to set the length of a sharable memory object.

4766 The meaning of *stat()* fields other than the size and protection information is undefined on
4767 implementations where memory objects are not implemented using regular files. When
4768 regular files are used, the times reflect when the implementation updated the file image of
4769 the data, not when a process updated the data in memory.

4770 The operations of *fdopen()*, *write()*, *read()*, and *lseek()* were made unspecified for objects
4771 opened with *shm_open()*, so that implementations that did not implement memory objects as
4772 regular files would not have to support the operation of these functions on shared memory
4773 objects.

4774 The behavior of memory objects with respect to *close()*, *dup()*, *dup2()*, *open()*, *close()*, *fork()*,
4775 *_exit()*, and the *exec* family of functions is the same as the behavior of the existing practice of
4776 the *mmap()* function.

4777 A memory object can still be referenced after a close. That is, any mappings made to the file
4778 are still in effect, and reads and writes that are made to those mappings are still valid and are
4779 shared with other processes that have the same mapping. Likewise, the memory object can
4780 still be used if any references remain after its name(s) have been deleted. Any references that
4781 remain after a close must not appear to the application as file descriptors.

4782 This is existing practice for *mmap()* and *close()*. In addition, there are already mappings
4783 present (text, data, stack) that do not have open file descriptors. The text mapping in
4784 particular is considered a reference to the file containing the text. The desire was to treat all
4785 mappings by the process uniformly. Also, many modern implementations use *mmap()* to
4786 implement shared libraries, and it would not be desirable to keep file descriptors for each of
4787 the many libraries an application can use. It was felt there were many other existing
4788 programs that used this behavior to free a file descriptor, and thus IEEE Std 1003.1-200x
4789 could not forbid it and still claim to be using existing practice.

4790 For implementations that implement memory objects using memory only, memory objects
4791 will retain the memory allocated to the file after the last close and will use that same memory
4792 on the next open. Note that closing the memory object is not the same as deleting the name,
4793 since the memory object is still defined in the memory object name space.

4794 The locks of *fcntl()* do not block any read or write operation, including read or write access to
4795 shared memory or mapped files. In addition, implementations that only support shared
4796 memory objects should not be required to implement record locks. The reference to *fcntl()* is
4797 added to make this point explicitly. The other *fcntl()* commands are useful with shared

4798 memory objects.

4799 The size of pages that mapping hardware may be able to support may be a configurable
4800 value, or it may change based on hardware implementations. The addition of the
4801 `_SC_PAGESIZE` parameter to the `sysconf()` function is provided for determining the mapping
4802 page size at runtime.

4803 **Shared Memory Functions**

4804 Implementations may support the Shared Memory Objects option without supporting a general
4805 Memory Mapped Files option. Shared memory objects are named regions of storage that may be
4806 independent of the file system and can be mapped into the address space of one or more
4807 processes to allow them to share the associated memory.

4808 • Requirements

4809 Shared memory is used to share data among several processes, each potentially running at
4810 different priority levels, responding to different inputs, or performing separate tasks. Shared
4811 memory is not just simply providing common access to data, it is providing the fastest
4812 possible communication between the processes. With one memory write operation, a process
4813 can pass information to as many processes as have the memory region mapped.

4814 As a result, shared memory provides a mechanism that can be used for all other interprocess
4815 communications facilities. It may also be used by an application for implementing more
4816 sophisticated mechanisms than semaphores and message queues.

4817 The need for a shared memory interface is obvious for virtual memory systems, where the
4818 operating system is directly preventing processes from accessing each other's data. However,
4819 in unprotected systems, such as those found in some embedded controllers, a shared
4820 memory interface is needed to provide a portable mechanism to allocate a region of memory
4821 to be shared and then to communicate the address of that region to other processes.

4822 This, then, provides the minimum functionality that a shared memory interface must have in
4823 order to support realtime applications: to allocate and name an object to be mapped into
4824 memory for potential sharing (`open()` or `shm_open()`), and to make the memory object
4825 available within the address space of a process (`mmap()`). To complete the interface, a
4826 mechanism to release the claim of a process on a shared memory object (`munmap()`) is also
4827 needed, as well as a mechanism for deleting the name of a sharable object that was
4828 previously created (`unlink()` or `shm_unlink()`).

4829 After a mapping has been established, an implementation should not have to provide
4830 services to maintain that mapping. All memory writes into that area will appear immediately
4831 in the memory mapping of that region by any other processes.

4832 Thus, requirements include:

4833 — Support creation of sharable memory objects and the mapping of these objects into the
4834 address space of a process.

4835 — Sharable memory objects should be accessed by global names accessible from all
4836 processes.

4837 — Support the mapping of specific sections of physical address space (such as a memory
4838 mapped device) into the address space of a process. This should not be done by the
4839 process specifying the actual address, but again by an implementation-defined global
4840 name (such as a special device name) dedicated to this purpose.

4841 — Support the mapping of discrete portions of these memory objects.

- 4842 — Support for minimum hardware configurations that contain no physical media on which
4843 to store shared memory contents permanently.
- 4844 — The ability to preallocate the entire shared memory region so that minimum hardware
4845 configurations without virtual memory support can guarantee contiguous space.
- 4846 — The maximizing of performance by not requiring functionality that would require
4847 implementation interaction above creating the shared memory area and returning the
4848 mapping.
- 4849 Note that the above requirements do not preclude:
- 4850 — The sharable memory object from being implemented using actual files on an actual file
4851 system.
- 4852 — The global name that is accessible from all processes being restricted to a file system area
4853 that is dedicated to handling shared memory.
- 4854 — An implementation not providing implementation-defined global names for the purpose
4855 of physical address mapping.
- 4856 • Shared Memory Objects Usage
- 4857 If the Shared Memory Objects option is supported, a shared memory object may be created,
4858 or opened if it already exists, with the *shm_open()* function. If the shared memory object is
4859 created, it has a length of zero. The *truncate()* function can be used to set the size of the
4860 shared memory object after creation. The *shm_unlink()* function removes the name for a
4861 shared memory object created by *shm_open()*.
- 4862 • Shared Memory Overview
- 4863 The shared memory facility defined by IEEE Std 1003.1-200x usually results in memory
4864 locations being added to the address space of the process. The implementation returns the
4865 address of the new space to the application by means of a pointer. This works well in
4866 languages like C. However, in languages without pointer types it will not work. In the
4867 bindings for such a language, either a special COMMON section will need to be defined
4868 (which is unlikely), or the binding will have to allow existing structures to be mapped. The
4869 implementation will likely have to place restrictions on the size and alignment of such
4870 structures or will have to map a suitable region of the address space of the process into the
4871 memory object, and thus into other processes. These are issues for that particular language
4872 binding. For IEEE Std 1003.1-200x, however, the practice will not be forbidden, merely
4873 undefined.
- 4874 Two potentially different name spaces are used for naming objects that may be mapped into
4875 process address spaces. When the Memory Mapped Files option is supported, files may be
4876 accessed via *open()*. When the Shared Memory Objects option is supported, sharable
4877 memory objects that might not be files may be accessed via the *shm_open()* function. These
4878 options are not mutually-exclusive.
- 4879 Some implementations supporting the Shared Memory Objects option may choose to |
4880 implement the shared memory object name space as part of the file system name space. |
4881 There are several reasons for this: |
- 4882 — It allows applications to prevent name conflicts by use of the directory structure.
- 4883 — It uses an existing mechanism for accessing global objects and prevents the creation of a
4884 new mechanism for naming global objects.
- 4885 In such implementations, memory objects can be implemented using regular files, if that is
4886 what the implementation chooses. The *shm_open()* function can be implemented as an *open()*

4887 call in a fixed directory followed by a call to *fcntl()* to set *FD_CLOEXEC*. The *shm_unlink()*
4888 function can be implemented as an *unlink()* call.

4889 On the other hand, it is also expected that small embedded systems that support the Shared
4890 Memory Objects option may wish to implement shared memory without having any file
4891 systems present. In this case, the implementations may choose to use a simple string valued
4892 name space for shared memory regions. The *shm_open()* function permits either type of
4893 implementation.

4894 Some implementations have hardware that supports protection of mapped data from certain |
4895 classes of access and some do not. Systems that supply this functionality can support the |
4896 Memory Protection option. |

4897 Some implementations restrict size, alignment, and protections to be on *page*-size
4898 boundaries. If an implementation has no restrictions on size or alignment, it may specify a 1-
4899 byte page size. Applications on implementations that do support larger pages must be
4900 cognizant of the page size since this is the alignment and protection boundary.

4901 Simple embedded implementations may have a 1-byte page size and only support the Shared
4902 Memory Objects option. This provides simple shared memory between processes without
4903 requiring mapping hardware.

4904 IEEE Std 1003.1-200x specifically allows a memory object to remain referenced after a close
4905 because that is existing practice for the *mmap()* function.

4906 **Typed Memory Functions**

4907 Implementations may support the Typed Memory Objects option without supporting either the
4908 Shared Memory option or the Memory Mapped Files option. Typed memory objects are pools of
4909 specialized storage, different from the main memory resource normally used by a processor to
4910 hold code and data, that can be mapped into the address space of one or more processes.

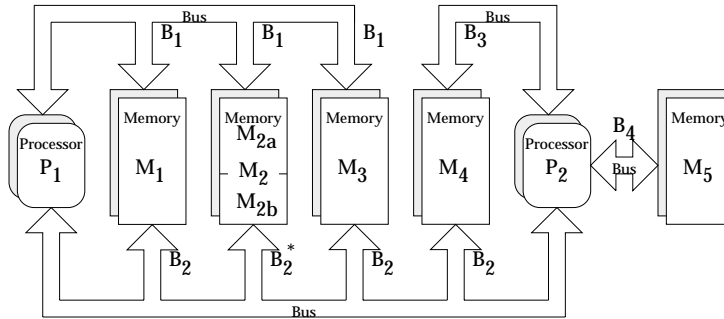
4911 • Model

4912 Realtime systems conforming to one of the POSIX.13 realtime profiles are expected (and
4913 desired) to be supported on systems with more than one type or pool of memory (for
4914 example, SRAM, DRAM, ROM, EPROM, EEPROM), where each type or pool of memory may
4915 be accessible by one or more processors via one or more busses (ports). Memory mapped
4916 files, shared memory objects, and the language-specific storage allocation operators (*malloc()*
4917 for the ISO C standard, *new* for ISO Ada) fail to provide application program interfaces
4918 versatile enough to allow applications to control their utilization of such diverse memory
4919 resources. The typed memory interfaces *posix_typed_mem_open()*, *posix_mem_offset()*,
4920 *posix_typed_mem_get_info()*, *mmap()*, and *munmap()* defined herein support the model of
4921 typed memory described below.

4922 For purposes of this model, a system comprises several processors (for example, P1 and P2),
4923 several physical memory pools (for example, M1, M2, M2a, M2b, M3, M4, and M5), and
4924 several busses or “ports” (for example, B1, B2, B3, and B4) interconnecting the various
4925 processors and memory pools in some system-specific way. Notice that some memory pools
4926 may be contained in others (for example, M2a and M2b are contained in M2).

4927 Figure B-1 (on page 3412) shows an example of such a model. In a system like this, an
4928 application should be able to perform the following operations:

4929



- * 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.

4930

Figure B-1 Example of a System with Typed Memory

4931

— Typed Memory Allocation

4932

4933

4934

4935

4936

4937

An application should be able to allocate memory dynamically from the desired pool using the desired bus, and map it into a process' address space. For example, processor 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 the process' address space, and subsequently deallocated when unmapped from all processes.

4938

— Using the Same Storage Region from Different Busses

4939

4940

4941

4942

4943

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 M2b both through ports B1 and B2.

4944

— Sharing Typed Memory Regions

4945

4946

4947

4948

4949

4950

4951

4952

4953

4954

4955

4956

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 busses 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.

4957

— Contiguous Allocation

4958

4959

4960

4961

The problem of finding the memory map of the different fragments of the memory pool that were mapped into logically contiguous addresses of a given process, can be solved by requesting contiguous allocation. For example, a process in P1 can allocate 10 Kbytes of physically contiguous memory from M3-B1, and obtain the offset (within pool M3) of

4962 this block of memory. Then, it can pass this offset (and the length) to a process in P2 using
 4963 some interprocess communication mechanism. The second process can map the same
 4964 block of memory by using the offset transferred and specifying M3-B2.

4965 — Unallocated Mapping

4966 Any subarea of a memory pool that is mapped to a process, either as the result of an
 4967 allocation request or an explicit mapping, is normally unavailable for allocation. Special
 4968 processes such as debuggers, however, may need to map large areas of a typed memory
 4969 pool, yet leave those areas available for allocation.

4970 Typed memory allocation and mapping has to coexist with storage allocation operators like
 4971 *malloc()*, but systems are free to choose how to implement this coexistence. For example, it
 4972 may be system configuration-dependent if all available system memory is made part of one
 4973 of the typed memory pools or if some part will be restricted to conventional allocation
 4974 operators. Equally system configuration-dependent may be the availability of operators like
 4975 *malloc()* to allocate storage from certain typed memory pools. It is not excluded to configure
 4976 a system such that a given named pool, P1, is in turn split into non-overlapping named
 4977 subpools. For example, M1-B1, M2-B1, and M3-B1 could also be accessed as one common
 4978 pool M123-B1. A call to *malloc()* on P1 could work on such a larger pool while full
 4979 optimization of memory usage by P1 would require typed memory allocation at the subpool
 4980 level.

4981 • Existing Practice

4982 OS-9 provides for the naming (numbering) and prioritization of memory types by a system
 4983 administrator. It then provides APIs to request memory allocation of typed (colored)
 4984 memory by number, and to generate a bus address from a mapped memory address
 4985 (translate). When requesting colored memory, the user can specify type 0 to signify allocation
 4986 from the first available type in priority order.

4987 HP-RT presents interfaces to map different kinds of storage regions that are visible through a
 4988 VME bus, although it does not provide allocation operations. It also provides functions to
 4989 perform address translation between VME addresses and virtual addresses. It represents a
 4990 VME-bus unique solution to the general problem.

4991 The PSOS approach is similar (that is, based on a pre-established mapping of bus address
 4992 ranges to specific memories) with a concept of segments and regions (regions dynamically
 4993 allocated from a heap which is a special segment). Therefore, PSOS does not fully address the
 4994 general allocation problem either. PSOS does not have a “process”-based model, but more of
 4995 a “thread”-only-based model of multi-tasking. So mapping to a process address space is not
 4996 an issue.

4997 QNX (a Canadian OS vendor specializing in realtime embedded systems on 80x86-based
 4998 processors) uses the System V approach of opening specially named devices (shared memory
 4999 segments) and using *mmap()* to then gain access from the process. They do not address
 5000 allocation directly, but once typed shared memory can be mapped, an “allocation manager”
 5001 process could be written to handle requests for allocation.

5002 The System V approach also included allocation, implemented by opening yet other special
 5003 “devices” which allocate, rather than appearing as a whole memory object.

5004 The Orkid realtime kernel interface definition has operations to manage memory “regions”
 5005 and “pools”, which are areas of memory that may reflect the differing physical nature of the
 5006 memory. Operations to allocate memory from these regions and pools are also provided. |

- 5007 • Requirements
- 5008 Existing practice in SVID-derived UNIX systems relies on functionality similar to *mmap()*
- 5009 and its related interfaces to achieve mapping and allocation of typed memory. However, the
- 5010 issue of sharing typed memory (allocated or mapped) and the complication of multiple ports
- 5011 are not addressed in any consistent way by existing UNIX system practice. Part of this
- 5012 functionality is existing practice in specialized realtime operating systems. In order to
- 5013 solidify the capabilities implied by the model above, the following requirements are imposed
- 5014 on the interface:
 - 5015 — Identification of Typed Memory Pools and Ports
 - 5016 All processes (running in all processors) in the system shall be able to identify a particular
 - 5017 (system configured) typed memory pool accessed through a particular (system
 - 5018 configured) port by a name. That name shall be a member of a name space common to all
 - 5019 these processes, but need not be the same name space as that containing ordinary
 - 5020 filenames. The association between memory pools/ports and corresponding names is
 - 5021 typically established when the system is configured. The “open” operation for typed
 - 5022 memory objects should be distinct from the *open()* function, for consistency with other
 - 5023 similar services, but implementable on top of *open()*. This implies that the handle for a
 - 5024 typed memory object will be a file descriptor.
 - 5025 — Allocation and Mapping of Typed Memory
 - 5026 Once a typed memory object has been identified by a process, it shall be possible to both
 - 5027 map user-selected subareas of that object into process address space and to map system-
 - 5028 selected (that is, dynamically allocated) subareas of that object, with user-specified
 - 5029 length, into process address space. It shall also be possible to determine the maximum
 - 5030 length of memory allocation that may be requested from a given typed memory object.
 - 5031 — Sharing Typed Memory
 - 5032 Two or more processes shall be able to share portions of typed memory, either user-
 - 5033 selected or dynamically allocated. This requirement applies also to dynamically allocated
 - 5034 regions of memory that are composed of several non-contiguous pieces.
 - 5035 — Contiguous Allocation
 - 5036 For dynamic allocation, it shall be the user’s option whether the system is required to
 - 5037 allocate a contiguous subarea within the typed memory object, or whether it is permitted
 - 5038 to allocate discontinuous fragments which appear contiguous in the process mapping.
 - 5039 Contiguous allocation simplifies the process of sharing allocated typed memory, while
 - 5040 discontinuous allocation allows for potentially better recovery of deallocated typed
 - 5041 memory.
 - 5042 — Accessing Typed Memory Through Different Ports
 - 5043 Once a subarea of a typed memory object has been mapped, it shall be possible to
 - 5044 determine the location and length corresponding to a user-selected portion of that object
 - 5045 within the memory pool. This location and length can then be used to remap that portion
 - 5046 of memory for access from another port. If the referenced portion of typed memory was
 - 5047 allocated discontinuously, the length thus determined may be shorter than anticipated,
 - 5048 and the user code shall adapt to the value returned.
 - 5049 — Deallocation
 - 5050 When a previously mapped subarea of typed memory is no longer mapped by any
 - 5051 process in the system—as a result of a call or calls to *munmap()*—that subarea shall
 - 5052 become potentially reusable for dynamic allocation; actual reuse of the subarea is a

- 5053 function of the dynamic typed memory allocation policy.
- 5054 — Unallocated Mapping
- 5055 It shall be possible to map user-selected subareas of a typed memory object without
5056 marking that subarea as unavailable for allocation. This option is not the default behavior,
5057 and shall require appropriate privilege.
- 5058 • Scenario
- 5059 The following scenario will serve to clarify the use of the typed memory interfaces.
- 5060 Process A running on P1 (see Figure B-1 (on page 3412)) wants to allocate some memory
5061 from memory pool M2, and it wants to share this portion of memory with process B running
5062 on P2. Since P2 only has access to the lower part of M2, both processes will use the memory
5063 pool named M2b which is the part of M2 that is accessible both from P1 and P2. The
5064 operations that both processes need to perform are shown below:
- 5065 — Allocating Typed Memory
- 5066 Process A calls *posix_typed_mem_open()* with the name */typed.m2b-b1* and a *tflag* of
5067 POSIX_TYPED_MEM_ALLOCATE to get a file descriptor usable for allocating from pool
5068 M2b accessed through port B1. It then calls *mmap()* with this file descriptor requesting a
5069 length of 4 096 bytes. The system allocates two discontinuous blocks of sizes 1 024 and
5070 3 072 bytes within M2b. The *mmap()* function returns a pointer to a 4 096 byte array in
5071 process A's logical address space, mapping the allocated blocks contiguously. Process A
5072 can then utilize the array, and store data in it.
- 5073 — Determining the Location of the Allocated Blocks
- 5074 Process A can determine the lengths and offsets (relative to M2b) of the two blocks
5075 allocated, by using the following procedure: First, process A calls *posix_mem_offset()*
5076 with the address of the first element of the array and length 4 096. Upon return, the offset and
5077 length (1 024 bytes) of the first block are returned. A second call to *posix_mem_offset()* is
5078 then made using the address of the first element of the array plus 1 024 (the length of the
5079 first block), and a new length of 4 096–1 024. If there were more fragments allocated, this
5080 procedure could have been continued within a loop until the offsets and lengths of all the
5081 blocks were obtained. Notice that this relatively complex procedure can be avoided if
5082 contiguous allocation is requested (by opening the typed memory object with the *tflag*
5083 POSIX_TYPED_MEM_ALLOCATE_CONTIG).
- 5084 — Sharing Data Across Processes
- 5085 Process A passes the two offset values and lengths obtained from the *posix_mem_offset()*
5086 calls to process B running on P2, via some form of interprocess communication. Process B
5087 can gain access to process A's data by calling *posix_typed_mem_open()* with the name
5088 */typed.m2b-b2* and a *tflag* of zero, then using two *mmap()* calls on the resulting file
5089 descriptor to map the two subareas of that typed memory object to its own address space.
- 5090 • Rationale for no *mem_alloc()* and *mem_free()*
- 5091 The standard developers had originally proposed a pair of new flags to *mmap()* which, when
5092 applied to a typed memory object descriptor, would cause *mmap()* to allocate dynamically
5093 from an unallocated and unmapped area of the typed memory object. Deallocation was
5094 similarly accomplished through the use of *munmap()*. This was rejected by the ballot group
5095 because it excessively complicated the (already rather complex) *mmap()* interface and
5096 introduced semantics useful only for typed memory, to a function which must also map
5097 shared memory and files. They felt that a memory allocator should be built on top of *mmap()*
5098 instead of being incorporated within the same interface, much as the ISO C standard libraries

5099 build *malloc()* on top of the virtual memory mapping functions *brk()* and *sbrk()*. This would
 5100 eliminate the complicated semantics involved with unmapping only part of an allocated
 5101 block of typed memory.

5102 To attempt to achieve ballot group consensus, typed memory allocation and deallocation was
 5103 first migrated from *mmap()* and *munmap()* to a pair of complementary functions modeled on
 5104 the ISO C standard *malloc()* and *free()*. The *mem_alloc()* function specified explicitly the
 5105 typed memory object (typed memory pool/access port) from which allocation takes place,
 5106 unlike *malloc()* where the memory pool and port are unspecified. The *mem_free()* function
 5107 handled deallocation. These new semantics still met all of the requirements detailed above
 5108 without modifying the behavior of *mmap()* except to allow it to map specified areas of typed
 5109 memory objects. An implementation would have been free to implement *mem_alloc()* and
 5110 *mem_free()* over *mmap()*, through *mmap()*, or independently but cooperating with *mmap()*.

5111 The ballot group was queried to see if this was an acceptable alternative, and while there was
 5112 some agreement that it achieved the goal of removing the complicated semantics of
 5113 allocation from the *mmap()* interface, several balloters realized that it just created two
 5114 additional functions that behaved, in great part, like *mmap()*. These balloters proposed an
 5115 alternative which has been implemented here in place of a separate *mem_alloc()* and
 5116 *mem_free()*. This alternative is based on four specific suggestions:

- 5117 1. The *posix_typed_mem_open()* function should provide a flag which specifies “allocate
 5118 on *mmap()*” (otherwise, *mmap()* just maps the underlying object). This allows things
 5119 roughly similar to */dev/zero* versus */dev/swap*. Two such flags have been implemented,
 5120 one of which forces contiguous allocation.
- 5121 2. The *posix_mem_offset()* function is acceptable because it can be applied usefully to
 5122 mapped objects in general. It should return the file descriptor of the underlying object.
- 5123 3. The *mem_get_info()* function in an earlier draft should be renamed
 5124 *posix_typed_mem_get_info()* because it is not generally applicable to memory objects. It
 5125 should probably return the file descriptor’s allocation attribute. We have implemented
 5126 the renaming of the function, but reject having it return a piece of information which is
 5127 readily known by an application without this function. Its whole purpose is to query
 5128 the typed memory object for attributes that are not user-specified, but determined by
 5129 the implementation.
- 5130 4. There should be no separate *mem_alloc()* or *mem_free()* functions. Instead, using
 5131 *mmap()* on a typed memory object opened with an “allocate on *mmap()*” flag should be
 5132 used to force allocation. These are precisely the semantics defined in the current draft.

5133 • Rationale for no Typed Memory Access Management

5134 The working group had originally defined an additional interface (and an additional kind of
 5135 object: typed memory master) to establish and dissolve mappings to typed memory on
 5136 behalf of devices or processors which were independent of the operating system and had no
 5137 inherent capability to directly establish mappings on their own. This was to have provided
 5138 functionality similar to device driver interfaces such as *physio()* and their underlying bus-
 5139 specific interfaces (for example, *mballoc()*) which serve to set up and break down DMA
 5140 pathways, and derive mapped addresses for use by hardware devices and processor cards.

5141 The ballot group felt that this was beyond the scope of POSIX.1 and its amendments.
 5142 Furthermore, the removal of interrupt handling interfaces from a preceding amendment (the
 5143 IEEE Std 1003.1d-1999) during its balloting process renders these typed memory access
 5144 management interfaces an incomplete solution to portable device management from a user
 5145 process; it would be possible to initiate a device transfer to/from typed memory, but
 5146 impossible to handle the transfer-complete interrupt in a portable way.

5147 To achieve ballot group consensus, all references to typed memory access management
 5148 capabilities were removed. The concept of portable interfaces from a device driver to both
 5149 operating system and hardware is being addressed by the Uniform Driver Interface (UDI)
 5150 industry forum, with formal standardization deferred until proof of concept and industry-
 5151 wide acceptance and implementation.

5152 B.2.8.4 Process Scheduling

5153 IEEE PASC Interpretation 1003.1 #96 has been applied, adding the `pthread_setschedprio()` |
 5154 function. This was added since previously there was no way for a thread to lower its own |
 5155 priority without going to the tail of the threads list for its new priority. This capability is |
 5156 necessary to bound the duration of priority inversion encountered by a thread. |

5157 The following portion of the rationale presents models, requirements, and standardization issues |
 5158 relevant to process scheduling; see also Section B.2.9.4 (on page 3456). |

5159 In an operating system supporting multiple concurrent processes, the system determines the
 5160 order in which processes execute to meet implementation-defined goals. For time-sharing |
 5161 systems, the goal is to enhance system throughput and promote fairness; the application is |
 5162 provided little or no control over this sequencing function. While this is acceptable and desirable |
 5163 behavior in a time-sharing system, it is inappropriate in a realtime system; realtime applications |
 5164 must specifically control the execution sequence of their concurrent processes in order to meet
 5165 externally defined response requirements.

5166 In IEEE Std 1003.1-200x, the control over process sequencing is provided using a concept of
 5167 scheduling policies. These policies, described in detail in this section, define the behavior of the
 5168 system whenever processor resources are to be allocated to competing processes. Only the
 5169 behavior of the policy is defined; conforming implementations are free to use any mechanism
 5170 desired to achieve the described behavior.

5171 • Models

5172 In an operating system supporting multiple concurrent processes, the system determines the
 5173 order in which processes execute and might force long-running processes to yield to other
 5174 processes at certain intervals. Typically, the scheduling code is executed whenever an event
 5175 occurs that might alter the process to be executed next.

5176 The simplest scheduling strategy is a “first-in, first-out” (FIFO) dispatcher. Whenever a
 5177 process becomes runnable, it is placed on the end of a ready list. The process at the front of
 5178 the ready list is executed until it exits or becomes blocked, at which point it is removed from
 5179 the list. This scheduling technique is also known as “run-to-completion” or “run-to-block”.

5180 A natural extension to this scheduling technique is the assignment of a “non-migrating
 5181 priority” to each process. This policy differs from strict FIFO scheduling in only one respect:
 5182 whenever a process becomes runnable, it is placed at the end of the list of processes runnable
 5183 at that priority level. When selecting a process to run, the system always selects the first
 5184 process from the highest priority queue with a runnable process. Thus, when a process
 5185 becomes unblocked, it will preempt a running process of lower priority without otherwise
 5186 altering the ready list. Further, if a process elects to alter its priority, it is removed from the
 5187 ready list and reinserted, using its new priority, according to the policy above.

5188 While the above policy might be considered unfriendly in a time-sharing environment in
 5189 which multiple users require more balanced resource allocation, it could be ideal in a
 5190 realtime environment for several reasons. The most important of these is that it is
 5191 deterministic: the highest-priority process is always run and, among processes of equal
 5192 priority, the process that has been runnable for the longest time is executed first. Because of
 5193 this determinism, cooperating processes can implement more complex scheduling simply by

- 5194 altering their priority. For instance, if processes at a single priority were to reschedule
5195 themselves at fixed time intervals, a time-slice policy would result.
- 5196 In a dedicated operating system in which all processes are well-behaved realtime
5197 applications, non-migrating priority scheduling is sufficient. However, many existing
5198 implementations provide for more complex scheduling policies.
- 5199 IEEE Std 1003.1-200x specifies a linear scheduling model. In this model, every process in the
5200 system has a priority. The system scheduler always dispatches a process that has the highest
5201 (generally the most time-critical) priority among all runnable processes in the system. As
5202 long as there is only one such process, the dispatching policy is trivial. When multiple
5203 processes of equal priority are eligible to run, they are ordered according to a strict run-to-
5204 completion (FIFO) policy.
- 5205 The priority is represented as a positive integer and is inherited from the parent process. For
5206 processes running under a fixed priority scheduling policy, the priority is never altered
5207 except by an explicit function call.
- 5208 It was determined arbitrarily that larger integers correspond to “higher priorities”.
- 5209 Certain implementations might impose restrictions on the priority ranges to which processes
5210 can be assigned. There also can be restrictions on the set of policies to which processes can be
5211 set.
- 5212 • Requirements
- 5213 Realtime processes require that scheduling be fast and deterministic, and that it guarantees
5214 to preempt lower priority processes.
- 5215 Thus, given the linear scheduling model, realtime processes require that they be run at a
5216 priority that is higher than other processes. Within this framework, realtime processes are
5217 free to yield execution resources to each other in a completely portable and implementation-
5218 defined manner.
- 5219 As there is a generally perceived requirement for processes at the same priority level to share
5220 processor resources more equitably, provisions are made by providing a scheduling policy
5221 (that is, SCHED_RR) intended to provide a timeslice-like facility.
- 5222 **Note:** The following topics assume that low numeric priority implies low scheduling criticality
5223 and *vice versa*.
- 5224 • Rationale for New Interface
- 5225 Realtime applications need to be able to determine when processes will run in relation to
5226 each other. It must be possible to guarantee that a critical process will run whenever it is
5227 runnable; that is, whenever it wants to for as long as it needs. SCHED_FIFO satisfies this
5228 requirement. Additionally, SCHED_RR was defined to meet a realtime requirement for a
5229 well-defined time-sharing policy for processes at the same priority.
- 5230 It would be possible to use the BSD *setpriority()* and *getpriority()* functions by redefining the
5231 meaning of the “nice” parameter according to the scheduling policy currently in use by the
5232 process. The System V *nice()* interface was felt to be undesirable for realtime because it
5233 specifies an adjustment to the “nice” value, rather than setting it to an explicit value.
5234 Realtime applications will usually want to set priority to an explicit value. Also, System V
5235 *nice()* does not allow for changing the priority of another process.
- 5236 With the POSIX.1b interfaces, the traditional “nice” value does not affect the SCHED_FIFO
5237 or SCHED_RR scheduling policies. If a “nice” value is supported, it is implementation-
5238 defined whether it affects the SCHED_OTHER policy.

5239 An important aspect of IEEE Std 1003.1-200x is the explicit description of the queuing and
5240 preemption rules. It is critical, to achieve deterministic scheduling, that such rules be stated
5241 clearly in IEEE Std 1003.1-200x.

5242 IEEE Std 1003.1-200x does not address the interaction between priority and swapping. The
5243 issues involved with swapping and virtual memory paging are extremely implementation-
5244 defined and would be nearly impossible to standardize at this point. The proposed
5245 scheduling paradigm, however, fully describes the scheduling behavior of runnable
5246 processes, of which one criterion is that the working set be resident in memory. Assuming
5247 the existence of a portable interface for locking portions of a process in memory, paging
5248 behavior need not affect the scheduling of realtime processes.

5249 IEEE Std 1003.1-200x also does not address the priorities of “system” processes. In general,
5250 these processes should always execute in low-priority ranges to avoid conflict with other
5251 realtime processes. Implementations should document the priority ranges in which system
5252 processes run.

5253 The default scheduling policy is not defined. The effect of I/O interrupts and other system
5254 processing activities is not defined. The temporary lending of priority from one process to
5255 another (such as for the purposes of affecting freeing resources) by the system is not
5256 addressed. Preemption of resources is not addressed. Restrictions on the ability of a process
5257 to affect other processes beyond a certain level (influence levels) is not addressed.

5258 The rationale used to justify the simple time-quantum scheduler is that it is common practice
5259 to depend upon this type of scheduling to assure “fair” distribution of processor resources
5260 among portions of the application that must interoperate in a serial fashion. Note that
5261 IEEE Std 1003.1-200x is silent with respect to the setting of this time quantum, or whether it is
5262 a system-wide value or a per-process value, although it appears that the prevailing realtime
5263 practice is for it to be a system-wide value.

5264 In a system with N processes at a given priority, all processor-bound, in which the time
5265 quantum is equal for all processes at a specific priority level, the following assumptions are
5266 made of such a scheduling policy:

- 5267 1. A time quantum Q exists and the current process will own control of the processor for
5268 at least a duration of Q and will have the processor for a duration of Q .
- 5269 2. The N th process at that priority will control a processor within a duration of $(N-1) \times Q$.

5270 These assumptions are necessary to provide equal access to the processor and bounded
5271 response from the application.

5272 The assumptions hold for the described scheduling policy only if no system overhead, such
5273 as interrupt servicing, is present. If the interrupt servicing load is non-zero, then one of the
5274 two assumptions becomes fallacious, based upon how Q is measured by the system.

5275 If Q is measured by clock time, then the assumption that the process obtains a duration Q
5276 processor time is false if interrupt overhead exists. Indeed, a scenario can be constructed with
5277 N processes in which a single process undergoes complete processor starvation if a
5278 peripheral device, such as an analog-to-digital converter, generates significant interrupt
5279 activity periodically with a period of $N \times Q$.

5280 If Q is measured as actual processor time, then the assumption that the N th process runs in
5281 within the duration $(N-1) \times Q$ is false.

5282 It should be noted that SCHED_FIFO suffers from interrupt-based delay as well. However,
5283 for SCHED_FIFO, the implied response of the system is “as soon as possible”, so that the
5284 interrupt load for this case is a vendor selection and not a compliance issue.

5285 With this in mind, it is necessary either to complete the definition by including bounds on the
5286 interrupt load, or to modify the assumptions that can be made about the scheduling policy.

5287 Since the motivation of inclusion of the policy is common usage, and since current
5288 applications do not enjoy the luxury of bounded interrupt load, item (2) above is sufficient to
5289 express existing application needs and is less restrictive in the standard definition. No
5290 difference in interface is necessary.

5291 In an implementation in which the time quantum is equal for all processes at a specific
5292 priority, our assumptions can then be restated as:

5293 — A time quantum Q exists, and a processor-bound process will be rescheduled after a
5294 duration of, at most, Q . Time quantum Q may be defined in either wall clock time or
5295 execution time.

5296 — In general, the N th process of a priority level should wait no longer than $(N-1) \times Q$ time
5297 to execute, assuming no processes exist at higher priority levels.

5298 — No process should wait indefinitely.

5299 For implementations supporting per-process time quanta, these assumptions can be readily
5300 extended.

5301 **Sporadic Server Scheduling Policy**

5302 The sporadic server is a mechanism defined for scheduling aperiodic activities in time-critical
5303 realtime systems. This mechanism reserves a certain bounded amount of execution capacity for
5304 processing aperiodic events at a high priority level. Any aperiodic events that cannot be
5305 processed within the bounded amount of execution capacity are executed in the background at a
5306 low priority level. Thus, a certain amount of execution capacity can be guaranteed to be
5307 available for processing periodic tasks, even under burst conditions in the arrival of aperiodic
5308 processing requests (that is, a large number of requests in a short time interval). The sporadic
5309 server also simplifies the schedulability analysis of the realtime system, because it allows
5310 aperiodic processes or threads to be treated as if they were periodic. The sporadic server was
5311 first described by Sprunt, et al.

5312 The key concept of the sporadic server is to provide and limit a certain amount of computation
5313 capacity for processing aperiodic events at their assigned normal priority, during a time interval
5314 called the *replenishment period*. Once the entity controlled by the sporadic server mechanism is
5315 initialized with its period and execution-time budget attributes, it preserves its execution
5316 capacity until an aperiodic request arrives. The request will be serviced (if there are no higher
5317 priority activities pending) as long as there is execution capacity left. If the request is completed,
5318 the actual execution time used to service it is subtracted from the capacity, and a replenishment
5319 of this amount of execution time is scheduled to happen one replenishment period after the
5320 arrival of the aperiodic request. If the request is not completed, because there is no execution
5321 capacity left, then the aperiodic process or thread is assigned a lower background priority. For
5322 each portion of consumed execution capacity the execution time used is replenished after one
5323 replenishment period. At the time of replenishment, if the sporadic server was executing at a
5324 background priority level, its priority is elevated to the normal level. Other similar
5325 replenishment policies have been defined, but the one presented here represents a compromise
5326 between efficiency and implementation complexity.

5327 The interface that appears in this section defines a new scheduling policy for threads and
5328 processes that behaves according to the rules of the sporadic server mechanism. Scheduling
5329 attributes are defined and functions are provided to allow the user to set and get the parameters
5330 that control the scheduling behavior of this mechanism, namely the normal and low priority, the
5331 replenishment period, the maximum number of pending replenishment operations, and the

- 5332 initial execution-time budget.
- 5333 • Scheduling Aperiodic Activities
- 5334 Virtually all realtime applications are required to process aperiodic activities. In many cases,
5335 there are tight timing constraints that the response to the aperiodic events must meet. Usual
5336 timing requirements imposed on the response to these events are:
- 5337 — The effects of an aperiodic activity on the response time of lower priority activities must
5338 be controllable and predictable.
 - 5339 — The system must provide the fastest possible response time to aperiodic events.
 - 5340 — It must be possible to take advantage of all the available processing bandwidth not
5341 needed by time-critical activities to enhance average-case response times to aperiodic
5342 events.
- 5343 Traditional methods for scheduling aperiodic activities are background processing, polling
5344 tasks, and direct event execution:
- 5345 — Background processing consists of assigning a very low priority to the processing of
5346 aperiodic events. It utilizes all the available bandwidth in the system that has not been
5347 consumed by higher priority threads. However, it is very difficult, or impossible, to meet
5348 requirements on average-case response time, because the aperiodic entity has to wait for
5349 the execution of all other entities which have higher priority.
 - 5350 — Polling consists of creating a periodic process or thread for servicing aperiodic requests.
5351 At regular intervals, the polling entity is started and it services accumulated pending
5352 aperiodic requests. If no aperiodic requests are pending, the polling entity suspends itself
5353 until its next period. Polling allows the aperiodic requests to be processed at a higher
5354 priority level. However, worst and average-case response times of polling entities are a
5355 direct function of the polling period, and there is execution overhead for each polling
5356 period, even if no event has arrived. If the deadline of the aperiodic activity is short
5357 compared to the inter-arrival time, the polling frequency must be increased to guarantee
5358 meeting the deadline. For this case, the increase in frequency can dramatically reduce the
5359 efficiency of the system and, therefore, its capacity to meet all deadlines. Yet, polling
5360 represents a good way to handle a large class of practical problems because it preserves
5361 system predictability, and because the amortized overhead drops as load increases.
 - 5362 — Direct event execution consists of executing the aperiodic events at a high fixed-priority
5363 level. Typically, the aperiodic event is processed by an interrupt service routine as soon as
5364 it arrives. This technique provides predictable response times for aperiodic events, but
5365 makes the response times of all lower priority activities completely unpredictable under
5366 burst arrival conditions. Therefore, if the density of aperiodic event arrivals is
5367 unbounded, it may be a dangerous technique for time-critical systems. Yet, for those cases
5368 in which the physics of the system imposes a bound on the event arrival rate, it is
5369 probably the most efficient technique.
 - 5370 — The sporadic server scheduling algorithm combines the predictability of the polling
5371 approach with the short response times of the direct event execution. Thus, it allows
5372 systems to meet an important class of application requirements that cannot be met by
5373 using the traditional approaches. Multiple sporadic servers with different attributes can
5374 be applied to the scheduling of multiple classes of aperiodic events, each with different
5375 kinds of timing requirements, such as individual deadlines, average response times, and
5376 so on. It also has many other interesting applications for realtime, such as scheduling
5377 producer/consumer tasks in time-critical systems, limiting the effects of faults on the
5378 estimation of task execution-time requirements, and so on.

- 5379 • Existing Practice
- 5380 The sporadic server has been used in different kinds of applications, including military
5381 avionics, robot control systems, industrial automation systems, and so on. There are
5382 examples of many systems that cannot be successfully scheduled using the classic
5383 approaches, such as direct event execution, or polling, and are schedulable using a sporadic
5384 server scheduler. The sporadic server algorithm itself can successfully schedule all systems
5385 scheduled with direct event execution or polling.
- 5386 The sporadic server scheduling policy has been implemented as a commercial product in the
5387 run-time system of the Verdex Ada compiler. There are also many applications that have
5388 used a much less efficient application-level sporadic server. These real-time applications
5389 would benefit from a sporadic server scheduler implemented at the scheduler level.
- 5390 • Library-Level *versus* Kernel-Level Implementation
- 5391 The sporadic server interface described in this section requires the sporadic server policy to
5392 be implemented at the same level as the scheduler. This means that the process sporadic
5393 server shall be implemented at the kernel level and the thread sporadic server policy shall be
5394 implemented at the same level as the thread scheduler; that is, kernel or library level.
- 5395 In an earlier interface for the sporadic server, this mechanism was implementable at a
5396 different level than the scheduler. This feature allowed the implementer to choose between
5397 an efficient scheduler-level implementation, or a simpler user or library-level
5398 implementation. However, the working group considered that this interface made the use of
5399 sporadic servers more complex, and that library-level implementations would lack some of
5400 the important functionality of the sporadic server, namely the limitation of the actual
5401 execution time of aperiodic activities. The working group also felt that the interface
5402 described in this chapter does not preclude library-level implementations of threads intended
5403 to provide efficient low-overhead scheduling for those threads that are not scheduled under
5404 the sporadic server policy.
- 5405 • Range of Scheduling Priorities
- 5406 Each of the scheduling policies supported in IEEE Std 1003.1-200x has an associated range of
5407 priorities. The priority ranges for each policy might or might not overlap with the priority
5408 ranges of other policies. For time-critical realtime applications it is usual for periodic and
5409 aperiodic activities to be scheduled together in the same processor. Periodic activities will
5410 usually be scheduled using the SCHED_FIFO scheduling policy, while aperiodic activities
5411 may be scheduled using SCHED_SPORADIC. Since the application developer will require
5412 complete control over the relative priorities of these activities in order to meet his timing
5413 requirements, it would be desirable for the priority ranges of SCHED_FIFO and
5414 SCHED_SPORADIC to overlap completely. Therefore, although IEEE Std 1003.1-200x does
5415 not require any particular relationship between the different priority ranges, it is
5416 recommended that these two ranges should coincide.
- 5417 • Dynamically Setting the Sporadic Server Policy
- 5418 Several members of the working group requested that implementations should not be
5419 required to support dynamically setting the sporadic server scheduling policy for a thread.
5420 The reason is that this policy may have a high overhead for library-level implementations of
5421 threads, and if threads are allowed to dynamically set this policy, this overhead can be
5422 experienced even if the thread does not use that policy. By disallowing the dynamic setting of
5423 the sporadic server scheduling policy, these implementations can accomplish efficient
5424 scheduling for threads using other policies. If a strictly conforming application needs to use
5425 the sporadic server policy, and is therefore willing to pay the overhead, it must set this policy
5426 at the time of thread creation.

5427 • Limitation of the Number of Pending Replenishments

5428 The number of simultaneously pending replenishment operations must be limited for each
 5429 sporadic server for two reasons: an unlimited number of replenishment operations would
 5430 need an unlimited number of system resources to store all the pending replenishment
 5431 operations; on the other hand, in some implementations each replenishment operation will
 5432 represent a source of priority inversion (just for the duration of the replenishment operation)
 5433 and thus, the maximum amount of replenishments must be bounded to guarantee bounded
 5434 response times. The way in which the number of replenishments is bounded is by lowering
 5435 the priority of the sporadic server to *sched_ss_low_priority* when the number of pending
 5436 replenishments has reached its limit. In this way, no new replenishments are scheduled until
 5437 the number of pending replenishments decreases.

5438 In the sporadic server scheduling policy defined in IEEE Std 1003.1-200x, the application can
 5439 specify the maximum number of pending replenishment operations for a single sporadic
 5440 server, by setting the value of the *sched_ss_max_repl* scheduling parameter. This value must
 5441 be between one and {SS_REPL_MAX}, which is a maximum limit imposed by the
 5442 implementation. The limit {SS_REPL_MAX} must be greater than or equal to
 5443 {_POSIX_SS_REPL_MAX}, which is defined to be four in IEEE Std 1003.1-200x. The minimum
 5444 limit of four was chosen so that an application can at least guarantee that four different
 5445 aperiodic events can be processed during each interval of length equal to the replenishment
 5446 period.

5447 *B.2.8.5 Clocks and Timers*

5448 • Clocks

5449 IEEE Std 1003.1-200x and the ISO C standard both define functions for obtaining system
 5450 time. Implicit behind these functions is a mechanism for measuring passage of time. This
 5451 specification makes this mechanism explicit and calls it a clock. The *CLOCK_REALTIME*
 5452 clock required by IEEE Std 1003.1-200x is a higher resolution version of the clock that
 5453 maintains POSIX.1 system time. This is a “system-wide” clock, in that it is visible to all
 5454 processes and, were it possible for multiple processes to all read the clock at the same time,
 5455 they would see the same value.

5456 An extensible interface was defined, with the ability for implementations to define additional
 5457 clocks. This was done because of the observation that many realtime platforms support
 5458 multiple clocks, and it was desired to fit this model within the standard interface. But
 5459 implementation-defined clocks need not represent actual hardware devices, nor are they
 5460 necessarily system-wide.

5461 • Timers

5462 Two timer types are required for a system to support realtime applications:

5463 1. One-shot

5464 A one-shot timer is a timer that is armed with an initial expiration time, either relative
 5465 to the current time or at an absolute time (based on some timing base, such as time in
 5466 seconds and nanoseconds since the Epoch). The timer expires once and then is
 5467 disarmed. With the specified facilities, this is accomplished by setting the *it_value*
 5468 member of the *value* argument to the desired expiration time and the *it_interval* member
 5469 to zero.

5470 2. Periodic

5471 A periodic timer is a timer that is armed with an initial expiration time, again either
 5472 relative or absolute, and a repetition interval. When the initial expiration occurs, the

5473 timer is reloaded with the repetition interval and continues counting. With the
 5474 specified facilities, this is accomplished by setting the *it_value* member of the *value*
 5475 argument to the desired initial expiration time and the *it_interval* member to the desired
 5476 repetition interval.

5477 For both of these types of timers, the time of the initial timer expiration can be specified in
 5478 two ways:

- 5479 1. Relative (to the current time)
- 5480 2. Absolute

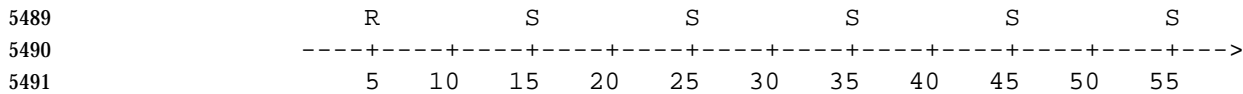
5481 • Examples of Using Realtime Timers

5482 In the diagrams below, *S* indicates a program schedule, *R* shows a schedule method request,
 5483 and *E* suggests an internal operating system event.

5484 — Periodic Timer: Data Logging

5485 During an experiment, it might be necessary to log realtime data periodically to an
 5486 internal buffer or to a mass storage device. With a periodic scheduling method, a logging
 5487 module can be started automatically at fixed time intervals to log the data.

5488 Program schedule is requested every 10 seconds.



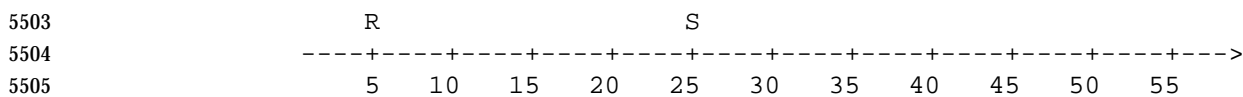
5492 [Time (in Seconds)]

5493 To achieve this type of scheduling using the specified facilities, one would allocate a per-
 5494 process timer based on clock ID `CLOCK_REALTIME`. Then the timer would be armed via
 5495 a call to `timer_settime()` with the `TIMER_ABSTIME` flag reset, and with an initial
 5496 expiration value and a repetition interval of 10 seconds.

5497 — One-shot Timer (Relative Time): Device Initialization

5498 In an emission test environment, large sample bags are used to capture the exhaust from
 5499 a vehicle. The exhaust is purged from these bags before each and every test. With a one-
 5500 shot timer, a module could initiate the purge function and then suspend itself for a
 5501 predetermined period of time while the sample bags are prepared.

5502 Program schedule requested 20 seconds after call is issued.



5506 [Time (in Seconds)]

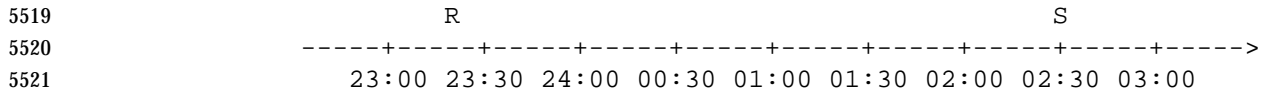
5507 To achieve this type of scheduling using the specified facilities, one would allocate a per-
 5508 process timer based on clock ID `CLOCK_REALTIME`. Then the timer would be armed via
 5509 a call to `timer_settime()` with the `TIMER_ABSTIME` flag reset, and with an initial
 5510 expiration value of 20 seconds and a repetition interval of zero.

5511 Note that if the program wishes merely to suspend itself for the specified interval, it
 5512 could more easily use `nanosleep()`.

5513 — One-shot Timer (Absolute Time): Data Transmission

5514 The results from an experiment are often moved to a different system within a network
 5515 for postprocessing or archiving. With an absolute one-shot timer, a module that moves
 5516 data from a test-cell computer to a host computer can be automatically scheduled on a
 5517 daily basis.

5518 Program schedule requested for 2:30 a.m.



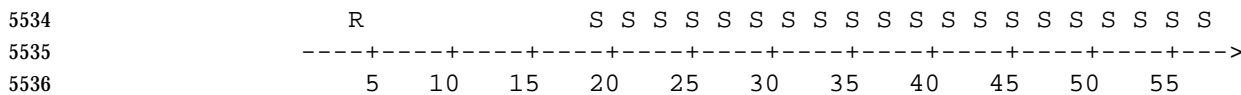
5522 [Time of Day]

5523 To achieve this type of scheduling using the specified facilities, one would allocate a per-
 5524 process timer based on clock ID CLOCK_REALTIME. Then the timer would be armed via
 5525 a call to *timer_settime()* with the *TIMER_ABSTIME* flag set, and an initial expiration value
 5526 equal to 2:30 a.m. of the next day.

5527 — Periodic Timer (Relative Time): Signal Stabilization

5528 Some measurement devices, such as emission analyzers, do not respond instantaneously
 5529 to an introduced sample. With a periodic timer with a relative initial expiration time, a
 5530 module that introduces a sample and records the average response could suspend itself
 5531 for a predetermined period of time while the signal is stabilized and then sample at a
 5532 fixed rate.

5533 Program schedule requested 15 seconds after call is issued and every 2 seconds thereafter.



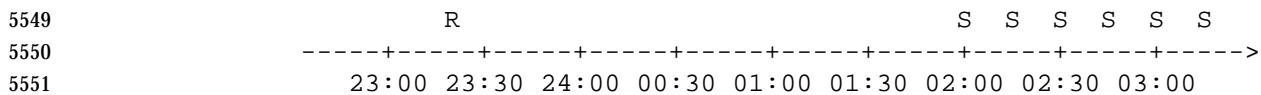
5537 [Time (in Seconds)]

5538 To achieve this type of scheduling using the specified facilities, one would allocate a per-
 5539 process timer based on clock ID CLOCK_REALTIME. Then the timer would be armed via
 5540 a call to *timer_settime()* with *TIMER_ABSTIME* flag reset, and with an initial expiration
 5541 value of 15 seconds and a repetition interval of 2 seconds.

5542 — Periodic Timer (Absolute Time): Work Shift-related Processing

5543 Resource utilization data is useful when time to perform experiments is being scheduled
 5544 at a facility. With a periodic timer with an absolute initial expiration time, a module can
 5545 be scheduled at the beginning of a work shift to gather resource utilization data
 5546 throughout the shift. This data can be used to allocate resources effectively to minimize
 5547 bottlenecks and delays and maximize facility throughput.

5548 Program schedule requested for 2:00 a.m. and every 15 minutes thereafter.



5552 [Time of Day]

5553 To achieve this type of scheduling using the specified facilities, one would allocate a per-
 5554 process timer based on clock ID CLOCK_REALTIME. Then the timer would be armed via
 5555 a call to *timer_settime()* with *TIMER_ABSTIME* flag set, and with an initial expiration
 5556 value equal to 2:00 a.m. and a repetition interval equal to 15 minutes.

5557 • Relationship of Timers to Clocks

5558 The relationship between clocks and timers armed with an absolute time is straightforward:
 5559 a timer expiration signal is requested when the associated clock reaches or exceeds the
 5560 specified time. The relationship between clocks and timers armed with a relative time (an
 5561 interval) is less obvious, but not unintuitive. In this case, a timer expiration signal is
 5562 requested when the specified interval, *as measured by the associated clock*, has passed. For the
 5563 required CLOCK_REALTIME clock, this allows timer expiration signals to be requested at
 5564 specified “wall clock” times (absolute), or when a specified interval of “realtime” has passed
 5565 (relative). For an implementation-defined clock—say, a process virtual time clock—timer
 5566 expirations could be requested when the process has used a specified total amount of virtual
 5567 time (absolute), or when it has used a specified *additional* amount of virtual time (relative).

5568 The interfaces also allow flexibility in the implementation of the functions. For example, an
 5569 implementation could convert all absolute times to intervals by subtracting the clock value at
 5570 the time of the call from the requested expiration time and “counting down” at the
 5571 supported resolution. Or it could convert all relative times to absolute expiration time by
 5572 adding in the clock value at the time of the call and comparing the clock value to the
 5573 expiration time at the supported resolution. Or it might even choose to maintain absolute
 5574 times as absolute and compare them to the clock value at the supported resolution for
 5575 absolute timers, and maintain relative times as intervals and count them down at the
 5576 resolution supported for relative timers. The choice will be driven by efficiency
 5577 considerations and the underlying hardware or software clock implementation.

5578 • Data Definitions for Clocks and Timers

5579 IEEE Std 1003.1-200x uses a time representation capable of supporting nanosecond resolution
 5580 timers for the following reasons:

- 5581 — To enable IEEE Std 1003.1-200x to represent those computer systems already using
 5582 nanosecond or submicrosecond resolution clocks.
- 5583 — To accommodate those per-process timers that might need nanoseconds to specify an
 5584 absolute value of system-wide clocks, even though the resolution of the per-process timer
 5585 may only be milliseconds, or *vice versa*.
- 5586 — Because the number of nanoseconds in a second can be represented in 32 bits.

5587 Time values are represented in the **timespec** structure. The *tv_sec* member is of type **time_t**
 5588 so that this member is compatible with time values used by POSIX.1 functions and the ISO C
 5589 standard. The *tv_nsec* member is a **signed long** in order to simplify and clarify code that
 5590 decrements or finds differences of time values. Note that because 1 billion (number of
 5591 nanoseconds per second) is less than half of the value representable by a signed 32-bit value,
 5592 it is always possible to add two valid fractional seconds represented as integral nanoseconds
 5593 without overflowing the signed 32-bit value.

5594 A maximum allowable resolution for the CLOCK_REALTIME clock of 20 ms (1/50 seconds)
 5595 was chosen to allow line frequency clocks in European countries to be conforming. 60 Hz
 5596 clocks in the U.S. will also be conforming, as will finer granularity clocks, although a Strictly
 5597 Conforming Application cannot assume a granularity of less than 20 ms (1/50 seconds).

5598 The minimum allowable maximum time allowed for the CLOCK_REALTIME clock and the
 5599 function *nanosleep()*, and timers created with *clock_id*=CLOCK_REALTIME, is determined by
 5600 the fact that the *tv_sec* member is of type **time_t**.

5601 IEEE Std 1003.1-200x specifies that timer expirations shall not be delivered early, nor shall
 5602 *nanosleep()* return early due to quantization error. IEEE Std 1003.1-200x discusses the various
 5603 implementations of *alarm()* in the rationale and states that implementations that do not

5604 allow alarm signals to occur early are the most appropriate, but refrained from mandating
5605 this behavior. Because of the importance of predictability to realtime applications,
5606 IEEE Std 1003.1-200x takes a stronger stance.

5607 The developers of IEEE Std 1003.1-200x considered using a time representation that differs
5608 from POSIX.1b in the second 32 bit of the 64-bit value. Whereas POSIX.1b defines this field
5609 as a fractional second in nanoseconds, the other methodology defines this as a binary fraction
5610 of one second, with the radix point assumed before the most significant bit.

5611 POSIX.1b is a software, source-level standard and most of the benefits of the alternate
5612 representation are enjoyed by hardware implementations of clocks and algorithms. It was
5613 felt that mandating this format for POSIX.1b clocks and timers would unnecessarily burden
5614 the application writer with writing, possibly non-portable, multiple precision arithmetic
5615 packages to perform conversion between binary fractions and integral units such as
5616 nanoseconds, milliseconds, and so on.

5617 **Rationale for the Monotonic Clock**

5618 For those applications that use time services to achieve realtime behavior, changing the value of
5619 the clock on which these services rely may cause erroneous timing behavior. For these
5620 applications, it is necessary to have a monotonic clock which cannot run backwards, and which
5621 has a maximum clock jump that is required to be documented by the implementation.
5622 Additionally, it is desirable (but not required by IEEE Std 1003.1-200x) that the monotonic clock
5623 increases its value uniformly. This clock should not be affected by changes to the system time;
5624 for example, to synchronize the clock with an external source or to account for leap seconds.
5625 Such changes would cause errors in the measurement of time intervals for those time services
5626 that use the absolute value of the clock.

5627 One could argue that by defining the behavior of time services when the value of a clock is
5628 changed, deterministic realtime behavior can be achieved. For example, one could specify that
5629 relative time services should be unaffected by changes in the value of a clock. However, there
5630 are time services that are based upon an absolute time, but that are essentially intended as
5631 relative time services. For example, *pthread_cond_timedwait()* uses an absolute time to allow it to
5632 wake up after the required interval despite spurious wakeups. Although sometimes the
5633 *pthread_cond_timedwait()* timeouts are absolute in nature, there are many occasions in which
5634 they are relative, and their absolute value is determined from the current time plus a relative
5635 time interval. In this latter case, if the clock changes while the thread is waiting, the wait interval
5636 will not be the expected length. If a *pthread_cond_timedwait()* function were created that would
5637 take a relative time, it would not solve the problem because to retain the intended “deadline” a
5638 thread would need to compensate for latency due to the spurious wakeup, and preemption
5639 between wakeup and the next wait.

5640 The solution is to create a new monotonic clock, whose value does not change except for the
5641 regular ticking of the clock, and use this clock for implementing the various relative timeouts
5642 that appear in the different POSIX interfaces, as well as allow *pthread_cond_timedwait()* to choose
5643 this new clock for its timeout. A new *clock_nanosleep()* function is created to allow an application
5644 to take advantage of this newly defined clock. Notice that the monotonic clock may be
5645 implemented using the same hardware clock as the system clock.

5646 Relative timeouts for *sigtimedwait()* and *aio_suspend()* have been redefined to use the monotonic
5647 clock, if present. The *alarm()* function has not been redefined, because the same effect but with
5648 better resolution can be achieved by creating a timer (for which the appropriate clock may be
5649 chosen).

5650 The *pthread_cond_timedwait()* function has been treated in a different way, compared to other
5651 functions with absolute timeouts, because it is used to wait for an event, and thus it may have a

5652 deadline, while the other timeouts are generally used as an error recovery mechanism, and for
5653 them the use of the monotonic clock is not so important. Since the desired timeout for the
5654 *pthread_cond_timedwait()* function may either be a relative interval, or an absolute time of day
5655 deadline, a new initialization attribute has been created for condition variables, to specify the
5656 clock that shall be used for measuring the timeout in a call to *pthread_cond_timedwait()*. In this
5657 way, if a relative timeout is desired, the monotonic clock will be used; if an absolute deadline is
5658 required instead, the `CLOCK_REALTIME` or another appropriate clock may be used. This
5659 capability has not been added to other functions with absolute timeouts because for those
5660 functions the expected use of the timeout is mostly to prevent errors, and not so often to meet
5661 precise deadlines. As a consequence, the complexity of adding this capability is not justified by
5662 its perceived application usage.

5663 The *nanosleep()* function has not been modified with the introduction of the monotonic clock.
5664 Instead, a new *clock_nanosleep()* function has been created, in which the desired clock may be
5665 specified in the function call.

5666 • History of Resolution Issues

5667 Due to the shift from relative to absolute timeouts in IEEE Std 1003.1d-1999, the amendments
5668 to the *sem_timedwait()*, *pthread_mutex_timedlock()*, *mq_timedreceive()*, and *mq_timedsend()*
5669 functions of that standard have been removed. Those amendments specified that
5670 `CLOCK_MONOTONIC` would be used for the (relative) timeouts if the Monotonic Clock
5671 option was supported.

5672 Having these functions continue to be tied solely to `CLOCK_MONOTONIC` would not
5673 work. Since the absolute value of a time value obtained from `CLOCK_MONOTONIC` is
5674 unspecified, under the absolute timeouts interface, applications would behave differently
5675 depending on whether the Monotonic Clock option was supported or not (because the
5676 absolute value of the clock would have different meanings in either case).

5677 Two options were considered:

- 5678 1. Leave the current behavior unchanged, which specifies the `CLOCK_REALTIME` clock
5679 for these (absolute) timeouts, to allow portability of applications between
5680 implementations supporting or not the Monotonic Clock option.
- 5681 2. Modify these functions in the way that *pthread_cond_timedwait()* was modified to allow
5682 a choice of clock, so that an application could use `CLOCK_REALTIME` when it is trying
5683 to achieve an absolute timeout and `CLOCK_MONOTONIC` when it is trying to achieve
5684 a relative timeout.

5685 It was decided that the features of `CLOCK_MONOTONIC` are not as critical to these
5686 functions as they are to *pthread_cond_timedwait()*. The *pthread_cond_timedwait()* function is
5687 given a relative timeout; the timeout may represent a deadline for an event. When these
5688 functions are given relative timeouts, the timeouts are typically for error recovery purposes
5689 and need not be so precise.

5690 Therefore, it was decided that these functions should be tied to `CLOCK_REALTIME` and not
5691 complicated by being given a choice of clock.

5692 **Execution Time Monitoring**

5693 • Introduction

5694 The main goals of the execution time monitoring facilities defined in this chapter are to
5695 measure the execution time of processes and threads and to allow an application to establish
5696 CPU time limits for these entities.

5697 The analysis phase of time-critical realtime systems often relies on the measurement of
5698 execution times of individual threads or processes to determine whether the timing
5699 requirements will be met. Also, performance analysis techniques for soft deadline realtime
5700 systems rely heavily on the determination of these execution times. The execution time
5701 monitoring functions provide application developers with the ability to measure these
5702 execution times online and open the possibility of dynamic execution-time analysis and
5703 system reconfiguration, if required.

5704 The second goal of allowing an application to establish execution time limits for individual
5705 processes or threads and detecting when they overrun allows program robustness to be
5706 increased by enabling online checking of the execution times.

5707 If errors are detected—possibly because of erroneous program constructs, the existence of
5708 errors in the analysis phase, or a burst of event arrivals—online detection and recovery is
5709 possible in a portable way. This feature can be extremely important for many time-critical
5710 applications. Other applications require trapping CPU-time errors as a normal way to exit an
5711 algorithm; for instance, some realtime artificial intelligence applications trigger a number of
5712 independent inference processes of varying accuracy and speed, limit how long they can run,
5713 and pick the best answer available when time runs out. In many periodic systems, overrun
5714 processes are simply restarted in the next resource period, after necessary end-of-period
5715 actions have been taken. This allows algorithms that are inherently data-dependent to be
5716 made predictable.

5717 The interface that appears in this chapter defines a new type of clock, the CPU-time clock,
5718 which measures execution time. Each process or thread can invoke the clock and timer
5719 functions defined in POSIX.1 to use them. Functions are also provided to access the CPU-
5720 time clock of other processes or threads to enable remote monitoring of these clocks.
5721 Monitoring of threads of other processes is not supported, since these threads are not visible
5722 from outside of their own process with the interfaces defined in POSIX.1.

5723 • Execution Time Monitoring Interface

5724 The clock and timer interface defined in POSIX.1 historically only defined one clock, which
5725 measures wall-clock time. The requirements for measuring execution time of processes and
5726 threads, and setting limits to their execution time by detecting when they overrun, can be
5727 accomplished with that interface if a new kind of clock is defined. These new clocks measure
5728 execution time, and one is associated with each process and with each thread. The clock
5729 functions currently defined in POSIX.1 can be used to read and set these CPU-time clocks,
5730 and timers can be created using these clocks as their timing base. These timers can then be
5731 used to send a signal when some specified execution time has been exceeded. The CPU-time
5732 clocks of each process or thread can be accessed by using the symbols
5733 CLOCK_PROCESS_CPUTIME_ID or CLOCK_THREAD_CPUTIME_ID.

5734 The clock and timer interface defined in POSIX.1 and extended with the new kind of CPU-
5735 time clock would only allow processes or threads to access their own CPU-time clocks.
5736 However, many realtime systems require the possibility of monitoring the execution time of
5737 processes or threads from independent monitoring entities. In order to allow applications to
5738 construct independent monitoring entities that do not require cooperation from or
5739 modification of the monitored entities, two functions have been added: *clock_getcpu* and *clockid()*,

5740 for accessing CPU-time clocks of other processes, and *pthread_getcpuclockid()*, for accessing
 5741 CPU-time clocks of other threads. These functions return the clock identifier associated with
 5742 the process or thread specified in the call. These clock IDs can then be used in the rest of the
 5743 clock function calls.

5744 The clocks accessed through these functions could also be used as a timing base for the
 5745 creation of timers, thereby allowing independent monitoring entities to limit the CPU-time
 5746 consumed by other entities. However, this possibility would imply additional complexity
 5747 and overhead because of the need to maintain a timer queue for each process or thread, to
 5748 store the different expiration times associated with timers created by different processes or
 5749 threads. The working group decided this additional overhead was not justified by
 5750 application requirements. Therefore, creation of timers attached to the CPU-time clocks of
 5751 other processes or threads has been specified as implementation-defined.

5752 • Overhead Considerations

5753 The measurement of execution time may introduce additional overhead in the thread
 5754 scheduling, because of the need to keep track of the time consumed by each of these entities.
 5755 In library-level implementations of threads, the efficiency of scheduling could be somehow
 5756 compromised because of the need to make a kernel call, at each context switch, to read the
 5757 process CPU-time clock. Consequently, a thread creation attribute called *cpu-clock-*
 5758 *requirement* was defined, to allow threads to disconnect their respective CPU-time clocks.
 5759 However, the Ballot Group considered that this attribute itself introduced some overhead,
 5760 and that in current implementations it was not worth the effort. Therefore, the attribute was
 5761 deleted, and thus thread CPU-time clocks are required for all threads if the Thread CPU-Time
 5762 Clocks option is supported.

5763 • Accuracy of CPU-time Clocks

5764 The mechanism used to measure the execution time of processes and threads is specified in
 5765 IEEE Std 1003.1-200x as implementation-defined. The reason for this is that both the
 5766 underlying hardware and the implementation architecture have a very strong influence on
 5767 the accuracy achievable for measuring CPU time. For some implementations, the
 5768 specification of strict accuracy requirements would represent very large overheads, or even
 5769 the impossibility of being implemented.

5770 Since the mechanism for measuring execution time is implementation-defined, realtime
 5771 applications will be able to take advantage of accurate implementations using a portable
 5772 interface. Of course, strictly conforming applications cannot rely on any particular degree of
 5773 accuracy, in the same way as they cannot rely on a very accurate measurement of wall clock
 5774 time. There will always exist applications whose accuracy or efficiency requirements on the
 5775 implementation are more rigid than the values defined in IEEE Std 1003.1-200x or any other
 5776 standard.

5777 In any case, there is a minimum set of characteristics that realtime applications would expect
 5778 from most implementations. One such characteristic is that the sum of all the execution times
 5779 of all the threads in a process equals the process execution time, when no CPU-time clocks
 5780 are disabled. This need not always be the case because implementations may differ in how
 5781 they account for time during context switches. Another characteristic is that the sum of the
 5782 execution times of all processes in a system equals the number of processors, multiplied by
 5783 the elapsed time, assuming that no processor is idle during that elapsed time. However, in |
 5784 some implementations it might not be possible to relate CPU-time to elapsed time. For |
 5785 example, in a heterogeneous multi-processor system in which each processor runs at a |
 5786 different speed, an implementation may choose to define each “second” of CPU-time to be a
 5787 certain number of “cycles” that a CPU has executed.

- 5788 • Existing Practice
- 5789 Measuring and limiting the execution time of each concurrent activity are common features
5790 of most industrial implementations of realtime systems. Almost all critical realtime systems
5791 are currently built upon a cyclic executive. With this approach, a regular timer interrupt kicks
5792 off the next sequence of computations. It also checks that the current sequence has
5793 completed. If it has not, then some error recovery action can be undertaken (or at least an
5794 overrun is avoided). Current software engineering principles and the increasing complexity
5795 of software are driving application developers to implement these systems on multi-
5796 threaded or multi-process operating systems. Therefore, if a POSIX operating system is to be
5797 used for this type of application, then it must offer the same level of protection.
- 5798 Execution time clocks are also common in most UNIX implementations, although these
5799 clocks usually have requirements different from those of realtime applications. The POSIX.1
5800 *times()* function supports the measurement of the execution time of the calling process, and
5801 its terminated child processes. This execution time is measured in clock ticks and is supplied
5802 as two different values with the user and system execution times, respectively. BSD supports
5803 the function *getrusage()*, which allows the calling process to get information about the
5804 resources used by itself and/or all of its terminated child processes. The resource usage
5805 includes user and system CPU time. Some UNIX systems have options to specify high
5806 resolution (up to one microsecond) CPU time clocks using the *times()* or the *getrusage()*
5807 functions.
- 5808 The *times()* and *getrusage()* interfaces do not meet important realtime requirements, such as
5809 the possibility of monitoring execution time from a different process or thread, or the
5810 possibility of detecting an execution time overrun. The latter requirement is supported in
5811 some UNIX implementations that are able to send a signal when the execution time of a
5812 process has exceeded some specified value. For example, BSD defines the functions
5813 *getitimer()* and *setitimer()*, which can operate either on a realtime clock (wall-clock), or on
5814 virtual-time or profile-time clocks which measure CPU time in two different ways. These
5815 functions do not support access to the execution time of other processes.
- 5816 IBM's MVS operating system supports per-process and per-thread execution time clocks. It
5817 also supports limiting the execution time of a given process.
- 5818 Given all this existing practice, the working group considered that the POSIX.1 clocks and
5819 timers interface was appropriate to meet most of the requirements that realtime applications
5820 have for execution time clocks. Functions were added to get the CPU time clock IDs, and to
5821 allow/disallow the thread CPU time clocks (in order to preserve the efficiency of some
5822 implementations of threads).
- 5823 • Clock Constants
- 5824 The definition of the manifest constants `CLOCK_PROCESS_CPUTIME_ID` and
5825 `CLOCK_THREAD_CPUTIME_ID` allows processes or threads, respectively, to access their
5826 own execution-time clocks. However, given a process or thread, access to its own execution-
5827 time clock is also possible if the clock ID of this clock is obtained through a call to
5828 *clock_getcpuclockid()* or *pthread_getcpuclockid()*. Therefore, these constants are not necessary
5829 and could be deleted to make the interface simpler. Their existence saves one system call in
5830 the first access to the CPU-time clock of each process or thread. The working group
5831 considered this issue and decided to leave the constants in IEEE Std 1003.1-200x because they
5832 are closer to the POSIX.1b use of clock identifiers.
- 5833 • Library Implementations of Threads
- 5834 In library implementations of threads, kernel entities and library threads can coexist. In this
5835 case, if the CPU-time clocks are supported, most of the clock and timer functions will need to

5836 have two implementations: one in the thread library, and one in the system calls library. The
5837 main difference between these two implementations is that the thread library
5838 implementation will have to deal with clocks and timers that reside in the thread space,
5839 while the kernel implementation will operate on timers and clocks that reside in kernel space.
5840 In the library implementation, if the clock ID refers to a clock that resides in the kernel, a
5841 kernel call will have to be made. The correct version of the function can be chosen by
5842 specifying the appropriate order for the libraries during the link process.

5843 • History of Resolution Issues: Deletion of the *enable* Attribute

5844 In the draft corresponding to the first balloting round, CPU-time clocks had an attribute
5845 called *enable*. This attribute was introduced by the working group to allow implementations
5846 to avoid the overhead of measuring execution time for those processes or threads for which
5847 this measurement was not required. However, the *enable* attribute got several ballot
5848 objections. The main reason was that processes are already required to measure execution
5849 time by the POSIX.1 *times()* function. Consequently, the *enable* attribute was considered
5850 unnecessary, and was deleted from the draft.

5851 Rationale Relating to Timeouts

5852 • Requirements for Timeouts

5853 Realtime systems which must operate reliably over extended periods without human
5854 intervention are characteristic in embedded applications such as avionics, machine control,
5855 and space exploration, as well as more mundane applications such as cable TV, security
5856 systems, and plant automation. A multi-tasking paradigm, in which many independent
5857 and/or cooperating software functions relinquish the processor(s) while waiting for a
5858 specific stimulus, resource, condition, or operation completion, is very useful in producing
5859 well engineered programs for such systems. For such systems to be robust and fault-tolerant,
5860 expected occurrences that are unduly delayed or that never occur must be detected so that
5861 appropriate recovery actions may be taken. This is difficult if there is no way for a task to
5862 regain control of a processor once it has relinquished control (blocked) awaiting an
5863 occurrence which, perhaps because of corrupted code, hardware malfunction, or latent
5864 software bugs, will not happen when expected. Therefore, the common practice in realtime
5865 operating systems is to provide a capability to timeout such blocking services. Although
5866 there are several methods to achieve this already defined by POSIX, none are as reliable or
5867 efficient as initiating a timeout simultaneously with initiating a blocking service. This is
5868 especially critical in hard-realtime embedded systems because the processors typically have
5869 little time reserve, and allowed fault recovery times are measured in milliseconds rather than
5870 seconds.

5871 The working group largely agreed that such timeouts were necessary and ought to become
5872 part of IEEE Std 1003.1-200x, particularly vendors of realtime operating systems whose
5873 customers had already expressed a strong need for timeouts. There was some resistance to
5874 inclusion of timeouts in IEEE Std 1003.1-200x because the desired effect, fault tolerance,
5875 could, in theory, be achieved using existing facilities and alternative software designs, but
5876 there was no compelling evidence that realtime system designers would embrace such
5877 designs at the sacrifice of performance and/or simplicity.

5878 • Which Services should be Timed Out?

5879 Originally, the working group considered the prospect of providing timeouts on all blocking
5880 services, including those currently existing in POSIX.1, POSIX.1b, and POSIX.1c, and future
5881 interfaces to be defined by other working groups, as sort of a general policy. This was rather
5882 quickly rejected because of the scope of such a change, and the fact that many of those
5883 services would not normally be used in a realtime context. More traditional timesharing

5884 solutions to timeout would suffice for most of the POSIX.1 interfaces, while others had
 5885 asynchronous alternatives which, while more complex to utilize, would be adequate for
 5886 some realtime and all non-realtime applications.

5887 The list of potential candidates for timeouts was narrowed to the following for further
 5888 consideration:

5889 — POSIX.1b

5890 — *sem_wait()*

5891 — *mq_receive()*

5892 — *mq_send()*

5893 — *lio_listio()*

5894 — *aio_suspend()*

5895 — *sigwait()* (timeout already implemented by *sigtimedwait()*)

5896 — POSIX.1c

5897 — *pthread_mutex_lock()*

5898 — *pthread_join()*

5899 — *pthread_cond_wait()* (timeout already implemented by *pthread_cond_timedwait()*)

5900 — POSIX.1

5901 — *read()*

5902 — *write()*

5903 After further review by the working group, the *lio_listio()*, *read()*, and *write()* functions (all
 5904 forms of blocking synchronous I/O) were eliminated from the list because of the following:

5905 — Asynchronous alternatives exist

5906 — Timeouts can be implemented, albeit non-portably, in device drivers

5907 — A strong desire not to introduce modifications to POSIX.1 interfaces

5908 The working group ultimately rejected *pthread_join()* since both that interface and a timed
 5909 variant of that interface are non-minimal and may be implemented as a function. See below
 5910 for a library implementation of *pthread_join()*.

5911 Thus, there was a consensus among the working group members to add timeouts to 4 of the
 5912 remaining 5 functions (the timeout for *aio_suspend()* was ultimately added directly to
 5913 POSIX.1b, while the others were added by POSIX.1d). However, *pthread_mutex_lock()*
 5914 remained contentious.

5915 Many feel that *pthread_mutex_lock()* falls into the same class as the other functions; that is, it
 5916 is desirable to timeout a mutex lock because a mutex may fail to be unlocked due to errant or
 5917 corrupted code in a critical section (looping or branching outside of the unlock code), and
 5918 therefore is equally in need of a reliable, simple, and efficient timeout. In fact, since mutexes
 5919 are intended to guard small critical sections, most *pthread_mutex_lock()* calls would be
 5920 expected to obtain the lock without blocking nor utilizing any kernel service, even in
 5921 implementations of threads with global contention scope; the timeout alternative need only
 5922 be considered after it is determined that the thread must block.

5923 Those opposed to timing out mutexes feel that the very simplicity of the mutex is
 5924 compromised by adding a timeout semantic, and that to do so is senseless. They claim that if

5925 a timed mutex is really deemed useful by a particular application, then it can be constructed
 5926 from the facilities already in POSIX.1b and POSIX.1c. The following two C-language library
 5927 implementations of mutex locking with timeout represent the solutions offered (in both
 5928 implementations, the timeout parameter is specified as absolute time, not relative time as in
 5929 the proposed POSIX.1c interfaces).

5930 • Spinlock Implementation

```

5931 #include <pthread.h>
5932 #include <time.h>
5933 #include <errno.h>

5934 int pthread_mutex_timedlock(pthread_mutex_t *mutex,
5935                             const struct timespec *timeout)
5936 {
5937     struct timespec timenow;

5938     while (pthread_mutex_trylock(mutex) == EBUSY)
5939     {
5940         clock_gettime(CLOCK_REALTIME, &timenow);
5941         if (timespec_cmp(&timenow, timeout) >= 0)
5942         {
5943             return ETIMEDOUT;
5944         }
5945         pthread_yield();
5946     }
5947     return 0;
5948 }
```

5949 The Spinlock implementation is generally unsuitable for any application using priority-based
 5950 thread scheduling policies such as SCHED_FIFO or SCHED_RR, since the mutex could
 5951 currently be held by a thread of lower priority within the same allocation domain, but since
 5952 the waiting thread never blocks, only threads of equal or higher priority will ever run, and
 5953 the mutex cannot be unlocked. Setting priority inheritance or priority ceiling protocol on the
 5954 mutex does not solve this problem, since the priority of a mutex owning thread is only
 5955 boosted if higher priority threads are blocked waiting for the mutex; clearly not the case for
 5956 this spinlock.

5957 • Condition Wait Implementation

```

5958 #include <pthread.h>
5959 #include <time.h>
5960 #include <errno.h>

5961 struct timed_mutex
5962 {
5963     int locked;
5964     pthread_mutex_t mutex;
5965     pthread_cond_t cond;
5966 };
5967 typedef struct timed_mutex timed_mutex_t;

5968 int timed_mutex_lock(timed_mutex_t *tm,
5969                    const struct timespec *timeout)
5970 {
5971     int timedout=FALSE;
5972     int error_status;
```

```

5973     pthread_mutex_lock(&tm->mutex);
5974     while (tm->locked && !timedout)
5975     {
5976         if ((error_status=pthread_cond_timedwait(&tm->cond,
5977         &tm->mutex,
5978         timeout))!=0)
5979         {
5980             if (error_status==ETIMEDOUT) timedout = TRUE;
5981         }
5982     }
5983     if(timedout)
5984     {
5985         pthread_mutex_unlock(&tm->mutex);
5986         return ETIMEDOUT;
5987     }
5988     else
5989     {
5990         tm->locked = TRUE;
5991         pthread_mutex_unlock(&tm->mutex);
5992         return 0;
5993     }
5994 }
5995 void timed_mutex_unlock(timed_mutex_t *tm)
5996 {
5997     pthread_mutex_lock(&tm->mutex); / for case assignment not atomic /
5998     tm->locked = FALSE;
5999     pthread_mutex_unlock(&tm->mutex);
6000     pthread_cond_signal(&tm->cond);
6001 }

```

6002 The Condition Wait implementation effectively substitutes the *pthread_cond_timedwait()*
6003 function (which is currently timed out) for the desired *pthread_mutex_timedlock()*. Since waits
6004 on condition variables currently do not include protocols which avoid priority inversion, this
6005 method is generally unsuitable for realtime applications because it does not provide the same
6006 priority inversion protection as the untimed *pthread_mutex_lock()*. Also, for any given
6007 implementations of the current mutex and condition variable primitives, this library
6008 implementation has a performance cost at least 2.5 times that of the untimed
6009 *pthread_mutex_lock()* even in the case where the timed mutex is readily locked without
6010 blocking (the interfaces required for this case are shown in bold). Even in uniprocessors or
6011 where assignment is atomic, at least an additional *pthread_cond_signal()* is required.
6012 *pthread_mutex_timedlock()* could be implemented at effectively no performance penalty in
6013 this case because the timeout parameters need only be considered after it is determined that
6014 the mutex cannot be locked immediately.

6015 Thus it has not yet been shown that the full semantics of mutex locking with timeout can be
6016 efficiently and reliably achieved using existing interfaces. Even if the existence of an
6017 acceptable library implementation were proven, it is difficult to justify why the interface
6018 itself should not be made portable, especially considering approval for the other four
6019 timeouts.

6020 • Rationale for Library Implementation of *pthread_timedjoin()*

```

6021     Library implementation of pthread_timedjoin():
6022     /*
6023     * Construct a thread variety entirely from existing functions
6024     * with which a join can be done, allowing the join to time out.
6025     */
6026     #include <pthread.h>
6027     #include <time.h>
6028
6028     struct timed_thread {
6029         pthread_t t;
6030         pthread_mutex_t m;
6031         int exiting;
6032         pthread_cond_t exit_c;
6033         void *(*start_routine)(void *arg);
6034         void *arg;
6035         void *status;
6036     };
6037
6037     typedef struct timed_thread *timed_thread_t;
6038     static pthread_key_t timed_thread_key;
6039     static pthread_once_t timed_thread_once = PTHREAD_ONCE_INIT;
6040
6040     static void timed_thread_init()
6041     {
6042         pthread_key_create(&timed_thread_key, NULL);
6043     }
6044
6044     static void *timed_thread_start_routine(void *args)
6045
6045     /*
6046     * Routine to establish thread-specific data value and run the actual
6047     * thread start routine which was supplied to timed_thread_create().
6048     */
6049     {
6050         timed_thread_t tt = (timed_thread_t) args;
6051
6051         pthread_once(&timed_thread_once, timed_thread_init);
6052         pthread_setspecific(timed_thread_key, (void *)tt);
6053         timed_thread_exit((tt->start_routine)(tt->arg));
6054     }
6055
6055     int timed_thread_create(timed_thread_t ttp, const pthread_attr_t *attr,
6056         void *(*start_routine)(void *), void *arg)
6057
6057     /*
6058     * Allocate a thread which can be used with timed_thread_join().
6059     */
6060     {
6061         timed_thread_t tt;
6062         int result;
6063
6063         tt = (timed_thread_t) malloc(sizeof(struct timed_thread));
6064         pthread_mutex_init(&tt->m, NULL);
6065         tt->exiting = FALSE;
6066         pthread_cond_init(&tt->exit_c, NULL);
6067         tt->start_routine = start_routine;

```

```

6068         tt->arg = arg;
6069         tt->status = NULL;

6070         if ((result = pthread_create(&tt->t, attr,
6071             timed_thread_start_routine, (void *)tt)) != 0) {
6072             free(tt);
6073             return result;
6074         }

6075         pthread_detach(tt->t);
6076         ttp = tt;
6077         return 0;
6078     }

6079     int timed_thread_join(timed_thread_t tt,
6080         struct timespec *timeout,
6081         void **status)
6082     {
6083         int result;

6084         pthread_mutex_lock(&tt->m);
6085         result = 0;
6086         /*
6087          * Wait until the thread announces that it is exiting,
6088          * or until timeout.
6089          */
6090         while (result == 0 && ! tt->exiting) {
6091             result = pthread_cond_timedwait(&tt->exit_c, &tt->m, timeout);
6092         }
6093         pthread_mutex_unlock(&tt->m);
6094         if (result == 0 && tt->exiting) {
6095             *status = tt->status;
6096             free((void *)tt);
6097             return result;
6098         }
6099         return result;
6100     }

6101     void timed_thread_exit(void *status)
6102     {
6103         timed_thread_t tt;
6104         void *specific;

6105         if ((specific=pthread_getspecific(timed_thread_key)) == NULL){
6106             /*
6107              * Handle cases which won't happen with correct usage.
6108              */
6109             pthread_exit( NULL);
6110         }
6111         tt = (timed_thread_t) specific;
6112         pthread_mutex_lock(&tt->m);
6113         /*
6114          * Tell a joiner that we're exiting.
6115          */
6116         tt->status = status;

```

```

6117         tt->exiting = TRUE;
6118         pthread_cond_signal(&tt->exit_c);
6119         pthread_mutex_unlock(&tt->m);
6120         /*
6121          * Call pthread_exit() to call destructors and really
6122          * exit the thread.
6123          */
6124         pthread_exit(NULL);
6125     }

```

6126 The *pthread_join()* C-language example shown above demonstrates that it is possible, using
6127 existing pthread facilities, to construct a variety of thread which allows for joining such a
6128 thread, but which allows the join operation to time out. It does this by using a
6129 *pthread_cond_timedwait()* to wait for the thread to exit. A **timed_thread_t** descriptor structure
6130 is used to pass parameters from the creating thread to the created thread, and from the
6131 exiting thread to the joining thread. This implementation is roughly equivalent to what a
6132 normal *pthread_join()* implementation would do, with the single change being that
6133 *pthread_cond_timedwait()* is used in place of a simple *pthread_cond_wait()*.

6134 Since it is possible to implement such a facility entirely from existing pthread interfaces, and
6135 with roughly equal efficiency and complexity to an implementation which would be
6136 provided directly by a pthreads implementation, it was the consensus of the working group
6137 members that any *pthread_timedjoin()* facility would be unnecessary, and should not be
6138 provided.

6139 • Form of the Timeout Interfaces

6140 The working group considered a number of alternative ways to add timeouts to blocking
6141 services. At first, a system interface which would specify a one-shot or persistent timeout to
6142 be applied to subsequent blocking services invoked by the calling process or thread was
6143 considered because it allowed all blocking services to be timed out in a uniform manner with
6144 a single additional interface; this was rather quickly rejected because it could easily result in
6145 the wrong services being timed out.

6146 It was suggested that a timeout value might be specified as an attribute of the object
6147 (semaphore, mutex, message queue, and so on), but there was no consensus on this, either on
6148 a case-by-case basis or for all timeouts.

6149 Looking at the two existing timeouts for blocking services indicates that the working group
6150 members favor a separate interface for the timed version of a function. However,
6151 *pthread_cond_timedwait()* utilizes an absolute timeout value while *sigtimedwait()* uses a
6152 relative timeout value. The working group members agreed that relative timeout values are
6153 appropriate where the timeout mechanism's primary use was to deal with an unexpected or
6154 error situation, but they are inappropriate when the timeout must expire at a particular time,
6155 or before a specific deadline. For the timeouts being introduced in IEEE Std 1003.1-200x, the
6156 working group considered allowing both relative and absolute timeouts as is done with
6157 POSIX.1b timers, but ultimately favored the simpler absolute timeout form.

6158 An absolute time measure can be easily implemented on top of an interface that specifies
6159 relative time, by reading the clock, calculating the difference between the current time and
6160 the desired wake-up time, and issuing a relative timeout call. But there is a race condition
6161 with this approach because the thread could be preempted after reading the clock, but before
6162 making the timed out call; in this case, the thread would be awakened later than it should
6163 and, thus, if the wake up time represented a deadline, it would miss it.

6164 There is also a race condition when trying to build a relative timeout on top of an interface
6165 that specifies absolute timeouts. In this case, we would have to read the clock to calculate the
6166 absolute wake-up time as the sum of the current time plus the relative timeout interval. In
6167 this case, if the thread is preempted after reading the clock but before making the timed out
6168 call, the thread would be awakened earlier than desired.

6169 But the race condition with the absolute timeouts interface is not as bad as the one that
6170 happens with the relative timeout interface, because there are simple workarounds. For the
6171 absolute timeouts interface, if the timing requirement is a deadline, we can still meet this
6172 deadline because the thread woke up earlier than the deadline. If the timeout is just used as
6173 an error recovery mechanism, the precision of timing is not really important. If the timing
6174 requirement is that between actions A and B a minimum interval of time must elapse, we can
6175 safely use the absolute timeout interface by reading the clock after action A has been started.
6176 It could be argued that, since the call with the absolute timeout is atomic from the
6177 application point of view, it is not possible to read the clock after action A, if this action is
6178 part of the timed out call. But if we look at the nature of the calls for which we specify
6179 timeouts (locking a mutex, waiting for a semaphore, waiting for a message, or waiting until
6180 there is space in a message queue), the timeouts that an application would build on these
6181 actions would not be triggered by these actions themselves, but by some other external
6182 action. For example, if we want to wait for a message to arrive to a message queue, and wait
6183 for at least 20 milliseconds, this time interval would start to be counted from some event that
6184 would trigger both the action that produces the message, as well as the action that waits for
6185 the message to arrive, and not by the wait-for-message operation itself. In this case, we could
6186 use the workaround proposed above.

6187 For these reasons, the absolute timeout is preferred over the relative timeout interface.

6188 **B.2.9 Threads**

6189 Threads will normally be more expensive than subroutines (or functions, routines, and so on) if
6190 specialized hardware support is not provided. Nevertheless, threads should be sufficiently
6191 efficient to encourage their use as a medium to fine-grained structuring mechanism for
6192 parallelism in an application. Structuring an application using threads then allows it to take
6193 immediate advantage of any underlying parallelism available in the host environment. This
6194 means implementors are encouraged to optimize for fast execution at the possible expense of
6195 efficient utilization of storage. For example, a common thread creation technique is to cache
6196 appropriate thread data structures. That is, rather than releasing system resources, the
6197 implementation retains these resources and reuses them when the program next asks to create a
6198 new thread. If this reuse of thread resources is to be possible, there has to be very little unique
6199 state associated with each thread, because any such state has to be reset when the thread is
6200 reused.

6201 **Thread Creation Attributes**

6202 Attributes objects are provided for threads, mutexes, and condition variables as a mechanism to
6203 support probable future standardization in these areas without requiring that the interface itself
6204 be changed. Attributes objects provide clean isolation of the configurable aspects of threads. For
6205 example, “stack size” is an important attribute of a thread, but it cannot be expressed portably.
6206 When porting a threaded program, stack sizes often need to be adjusted. The use of attributes
6207 objects can help by allowing the changes to be isolated in a single place, rather than being spread
6208 across every instance of thread creation.

6209 Attributes objects can be used to set up *classes* of threads with similar attributes; for example,
6210 “threads with large stacks and high priority” or “threads with minimal stacks”. These classes
6211 can be defined in a single place and then referenced wherever threads need to be created.

6212 Changes to “class” decisions become straightforward, and detailed analysis of each
6213 *pthread_create()* call is not required.

6214 The attributes objects are defined as opaque types as an aid to extensibility. If these objects had
6215 been specified as structures, adding new attributes would force recompilation of all multi-
6216 threaded programs when the attributes objects are extended; this might not be possible if
6217 different program components were supplied by different vendors.

6218 Additionally, opaque attributes objects present opportunities for improving performance.
6219 Argument validity can be checked once when attributes are set, rather than each time a thread is
6220 created. Implementations will often need to cache kernel objects that are expensive to create.
6221 Opaque attributes objects provide an efficient mechanism to detect when cached objects become
6222 invalid due to attribute changes.

6223 Because assignment is not necessarily defined on a given opaque type, implementation-
6224 dependent default values cannot be defined in a portable way. The solution to this problem is to
6225 allow attribute objects to be initialized dynamically by attributes object initialization functions,
6226 so that default values can be supplied automatically by the implementation.

6227 The following proposal was provided as a suggested alternative to the supplied attributes:

- 6228 1. Maintain the style of passing a parameter formed by the bitwise-inclusive OR of flags to
6229 the initialization routines (*pthread_create()*, *pthread_mutex_init()*, *pthread_cond_init()*). The
6230 parameter containing the flags should be an opaque type for extensibility. If no flags are
6231 set in the parameter, then the objects are created with default characteristics. An
6232 implementation may specify implementation-defined flag values and associated behavior.
- 6233 2. If further specialization of mutexes and condition variables is necessary, implementations
6234 may specify additional procedures that operate on the **pthread_mutex_t** and
6235 **pthread_cond_t** objects (instead of on attributes objects).

6236 The difficulties with this solution are:

- 6237 1. A bitmask is not opaque if bits have to be set into bit-vector attributes objects using
6238 explicitly-coded bitwise-inclusive OR operations. If the set of options exceeds an **int**,
6239 application programmers need to know the location of each bit. If bits are set or read by
6240 encapsulation (that is, *get*()* or *set*()* functions), then the bitmask is merely an
6241 implementation of attributes objects as currently defined and should not be exposed to the
6242 programmer.
- 6243 2. Many attributes are not Boolean or very small integral values. For example, scheduling
6244 policy may be placed in 3 bits or 4 bits, but priority requires 5 bits or more, thereby taking
6245 up at least 8 bits out of a possible 16 bits on machines with 16-bit integers. Because of this,
6246 the bitmask can only reasonably control whether particular attributes are set or not, and it
6247 cannot serve as the repository of the value itself. The value needs to be specified as a
6248 function parameter (which is non-extensible), or by setting a structure field (which is non-
6249 opaque), or by *get*()* and *set*()* functions (making the bitmask a redundant addition to the
6250 attributes objects).

6251 Stack size is defined as an optional attribute because the very notion of a stack is inherently
6252 machine-dependent. Some implementations may not be able to change the size of the stack, for
6253 example, and others may not need to because stack pages may be discontinuous and can be
6254 allocated and released on demand.

6255 The attribute mechanism has been designed in large measure for extensibility. Future extensions
6256 to the attribute mechanism or to any attributes object defined in IEEE Std 1003.1-200x has to be
6257 done with care so as not to affect binary-compatibility.

6258 Attribute objects, even if allocated by means of dynamic allocation functions such as *malloc()*,
6259 may have their size fixed at compile time. This means, for example, a *pthread_create()* in an
6260 implementation with extensions to the **pthread_attr_t** cannot look beyond the area that the
6261 binary application assumes is valid. This suggests that implementations should maintain a size
6262 field in the attributes object, as well as possibly version information, if extensions in different
6263 directions (possibly by different vendors) are to be accommodated.

6264 **Thread Implementation Models**

6265 There are various thread implementation models. At one end of the spectrum is the “library-
6266 thread model”. In such a model, the threads of a process are not visible to the operating system
6267 kernel, and the threads are not kernel scheduled entities. The process is the only kernel
6268 scheduled entity. The process is scheduled onto the processor by the kernel according to the
6269 scheduling attributes of the process. The threads are scheduled onto the single kernel scheduled
6270 entity (the process) by the runtime library according to the scheduling attributes of the threads.
6271 A problem with this model is that it constrains concurrency. Since there is only one kernel
6272 scheduled entity (namely, the process), only one thread per process can execute at a time. If the
6273 thread that is executing blocks on I/O, then the whole process blocks.

6274 At the other end of the spectrum is the “kernel-thread model”. In this model, all threads are
6275 visible to the operating system kernel. Thus, all threads are kernel scheduled entities, and all
6276 threads can concurrently execute. The threads are scheduled onto processors by the kernel
6277 according to the scheduling attributes of the threads. The drawback to this model is that the
6278 creation and management of the threads entails operating system calls, as opposed to subroutine
6279 calls, which makes kernel threads heavier weight than library threads.

6280 Hybrids of these two models are common. A hybrid model offers the speed of library threads
6281 and the concurrency of kernel threads. In hybrid models, a process has some (relatively small)
6282 number of kernel scheduled entities associated with it. It also has a potentially much larger
6283 number of library threads associated with it. Some library threads may be bound to kernel
6284 scheduled entities, while the other library threads are multiplexed onto the remaining kernel
6285 scheduled entities. There are two levels of thread scheduling:

- 6286 1. The runtime library manages the scheduling of (unbound) library threads onto kernel
6287 scheduled entities.
- 6288 2. The kernel manages the scheduling of kernel scheduled entities onto processors.

6289 For this reason, a hybrid model is referred to as a *two-level threads scheduling model*. In this model,
6290 the process can have multiple concurrently executing threads; specifically, it can have as many
6291 concurrently executing threads as it has kernel scheduled entities.

6292 **Thread-Specific Data**

6293 Many applications require that a certain amount of context be maintained on a per-thread basis
6294 across procedure calls. A common example is a multi-threaded library routine that allocates
6295 resources from a common pool and maintains an active resource list for each thread. The
6296 thread-specific data interface provided to meet these needs may be viewed as a two-dimensional
6297 array of values with keys serving as the row index and thread IDs as the column index (although
6298 the implementation need not work this way).

6299 • Models

6300 Three possible thread-specific data models were considered:

- 6301 1. No Explicit Support

6302 A standard thread-specific data interface is not strictly necessary to support
6303 applications that require per-thread context. One could, for example, provide a hash
6304 function that converted a `pthread_t` into an integer value that could then be used to
6305 index into a global array of per-thread data pointers. This hash function, in conjunction
6306 with `pthread_self()`, would be all the interface required to support a mechanism of this
6307 sort. Unfortunately, this technique is cumbersome. It can lead to duplicated code as
6308 each set of cooperating modules implements their own per-thread data management
6309 schemes.

6310 2. Single (`void *`) Pointer

6311 Another technique would be to provide a single word of per-thread storage and a pair
6312 of functions to fetch and store the value of this word. The word could then hold a
6313 pointer to a block of per-thread memory. The allocation, partitioning, and general use
6314 of this memory would be entirely up to the application. Although this method is not as
6315 problematic as technique 1, it suffers from interoperability problems. For example, all
6316 modules using the per-thread pointer would have to agree on a common usage
6317 protocol.

6318 3. Key/Value Mechanism

6319 This method associates an opaque key (for example, stored in a variable of type
6320 `pthread_key_t`) with each per-thread datum. These keys play the role of identifiers for
6321 per-thread data. This technique is the most generic and avoids the problems noted
6322 above, albeit at the cost of some complexity.

6323 The primary advantage of the third model is its information hiding properties. Modules
6324 using this model are free to create and use their own key(s) independent of all other such
6325 usage, whereas the other models require that all modules that use thread-specific context
6326 explicitly cooperate with all other such modules. The data-independence provided by the
6327 third model is worth the additional interface.

6328 • Requirements

6329 It is important that it be possible to implement the thread-specific data interface without the
6330 use of thread private memory. To do otherwise would increase the weight of each thread,
6331 thereby limiting the range of applications for which the threads interfaces provided by
6332 IEEE Std 1003.1-200x is appropriate.

6333 The values that one binds to the key via `pthread_setspecific()` may, in fact, be pointers to
6334 shared storage locations available to all threads. It is only the key/value bindings that are
6335 maintained on a per-thread basis, and these can be kept in any portion of the address space
6336 that is reserved for use by the calling thread (for example, on the stack). Thus, no per-thread
6337 MMU state is required to implement the interface. On the other hand, there is nothing in the
6338 interface specification to preclude the use of a per-thread MMU state if it is available (for
6339 example, the key values returned by `pthread_key_create()` could be thread private memory
6340 addresses).

6341 • Standardization Issues

6342 Thread-specific data is a requirement for a usable thread interface. The binding described in
6343 this section provides a portable thread-specific data mechanism for languages that do not
6344 directly support a thread-specific storage class. A binding to IEEE Std 1003.1-200x for a
6345 language that does include such a storage class need not provide this specific interface.

6346 If a language were to include the notion of thread-specific storage, it would be desirable (but
6347 *not* required) to provide an implementation of the pthreads thread-specific data interface
6348 based on the language feature. For example, assume that a compiler for a C-like language

6349 supports a *private* storage class that provides thread-specific storage. Something similar to
 6350 the following macros might be used to effect a compatible implementation:

```
6351     #define pthread_key_t                private void *
6352     #define pthread_key_create(key)      /* no-op */
6353     #define pthread_setspecific(key,value) (key)=(value)
6354     #define pthread_getspecific(key)     (key)
```

6355 **Note:** For the sake of clarity, this example ignores destructor functions. A correct implementation
 6356 would have to support them.

6357 Barriers

6358 • Background

6359 Barriers are typically used in parallel DO/FOR loops to ensure that all threads have reached
 6360 a particular stage in a parallel computation before allowing any to proceed to the next stage.
 6361 Highly efficient implementation is possible on machines which support a “Fetch and Add”
 6362 operation as described in the referenced Almasi and Gottlieb (1989).

6363 The use of return value PTHREAD_BARRIER_SERIAL_THREAD is shown in the following
 6364 example:

```
6365     if ( (status=pthread_barrier_wait(&barrier)) ==
6366         PTHREAD_BARRIER_SERIAL_THREAD) {
6367         ...serial section
6368     }
6369         else if (status != 0) {
6370         ...error processing
6371     }
6372     status=pthread_barrier_wait(&barrier);
6373     ...
```

6374 This behavior allows a serial section of code to be executed by one thread as soon as all
 6375 threads reach the first barrier. The second barrier prevents the other threads from proceeding
 6376 until the serial section being executed by the one thread has completed.

6377 Although barriers can be implemented with mutexes and condition variables, the referenced
 6378 Almasi and Gottlieb (1989) provides ample illustration that such implementations are
 6379 significantly less efficient than is possible. While the relative efficiency of barriers may well
 6380 vary by implementation, it is important that they be recognized in the IEEE Std 1003.1-200x
 6381 to facilitate application portability while providing the necessary freedom to implementors.

6382 • Lack of Timeout Feature

6383 Alternate versions of most blocking routines have been provided to support watchdog
 6384 timeouts. No alternate interface of this sort has been provided for barrier waits for the
 6385 following reasons:

- 6386 • Multiple threads may use different timeout values, some of which may be indefinite. It is
 6387 not clear which threads should break through the barrier with a timeout error if and when
 6388 these timeouts expire.
- 6389 • The barrier may become unusable once a thread breaks out of a *pthread_barrier_wait()*
 6390 with a timeout error. There is, in general, no way to guarantee the consistency of a
 6391 barrier’s internal data structures once a thread has timed out of a *pthread_barrier_wait()*.
 6392 Even the inclusion of a special barrier reinitialization function would not help much since
 6393 it is not clear how this function would affect the behavior of threads that reach the barrier

6394 between the original timeout and the call to the reinitialization function.

6395 **Spin Locks**

6396 • Background

6397 Spin locks represent an extremely low-level synchronization mechanism suitable primarily
6398 for use on shared memory multi-processors. It is typically an atomically modified Boolean
6399 value that is set to one when the lock is held and to zero when the lock is freed.

6400 When a caller requests a spin lock that is already held, it typically spins in a loop testing
6401 whether the lock has become available. Such spinning wastes processor cycles so the lock
6402 should only be held for short durations and not across sleep/block operations. Callers should
6403 unlock spin locks before calling sleep operations.

6404 Spin locks are available on a variety of systems. The functions included in
6405 IEEE Std 1003.1-200x are an attempt to standardize that existing practice.

6406 • Lack of Timeout Feature

6407 Alternate versions of most blocking routines have been provided to support watchdog
6408 timeouts. No alternate interface of this sort has been provided for spin locks for the following
6409 reasons:

6410 • It is impossible to determine appropriate timeout intervals for spin locks in a portable
6411 manner. The amount of time one can expect to spend spin-waiting is inversely
6412 proportional to the degree of parallelism provided by the system.

6413 It can vary from a few cycles when each competing thread is running on its own
6414 processor, to an indefinite amount of time when all threads are multiplexed on a single
6415 processor (which is why spin locking is not advisable on uniprocessors).

6416 • When used properly, the amount of time the calling thread spends waiting on a spin lock
6417 should be considerably less than the time required to set up a corresponding watchdog
6418 timer. Since the primary purpose of spin locks is to provide a low-overhead
6419 synchronization mechanism for multi-processors, the overhead of a timeout mechanism
6420 was deemed unacceptable.

6421 It was also suggested that an additional *count* argument be provided (on the
6422 *pthread_spin_lock()* call) in *lieu* of a true timeout so that a spin lock call could fail gracefully if
6423 it was unable to apply the lock after *count* attempts. This idea was rejected because it is not
6424 existing practice. Furthermore, the same effect can be obtained with *pthread_spin_trylock()*,
6425 as illustrated below:

```

6426         int n = MAX_SPIN;
6427         while ( --n >= 0 )
6428         {
6429             if ( !pthread_spin_try_lock(...) )
6430                 break;
6431         }
6432         if ( n >= 0 )
6433         {
6434             /* Successfully acquired the lock */
6435         }
6436         else
6437         {
6438             /* Unable to acquire the lock */
6439         }

```

6440 • *process-shared* Attribute

6441 The initialization functions associated with most POSIX synchronization objects (for
6442 example, mutexes, barriers, and read-write locks) take an attributes object with a *process-*
6443 *shared* attribute that specifies whether or not the object is to be shared across processes. In the
6444 draft corresponding to the first balloting round, two separate initialization functions are
6445 provided for spin locks, however: one for spin locks that were to be shared across processes
6446 (*spin_init()*), and one for locks that were only used by multiple threads within a single
6447 process (*pthread_spin_init()*). This was done so as to keep the overhead associated with spin
6448 waiting to an absolute minimum. However, the balloting group requested that, since the
6449 overhead associated to a bit check was small, spin locks should be consistent with the rest of
6450 the synchronization primitives, and thus the *process-shared* attribute was introduced for spin
6451 locks.

6452 • Spin Locks *versus* Mutexes

6453 It has been suggested that mutexes are an adequate synchronization mechanism and spin
6454 locks are not necessary. Locking mechanisms typically must trade off the processor resources
6455 consumed while setting up to block the thread and the processor resources consumed by the
6456 thread while it is blocked. Spin locks require very little resources to set up the blocking of a
6457 thread. Existing practice is to simply loop, repeating the atomic locking operation until the
6458 lock is available. While the resources consumed to set up blocking of the thread are low, the
6459 thread continues to consume processor resources while it is waiting.

6460 On the other hand, mutexes may be implemented such that the processor resources
6461 consumed to block the thread are large relative to a spin lock. After detecting that the mutex
6462 lock is not available, the thread must alter its scheduling state, add itself to a set of waiting
6463 threads, and, when the lock becomes available again, undo all of this before taking over
6464 ownership of the mutex. However, while a thread is blocked by a mutex, no processor
6465 resources are consumed.

6466 Therefore, spin locks and mutexes may be implemented to have different characteristics.
6467 Spin locks may have lower overall overhead for very short-term blocking, and mutexes may
6468 have lower overall overhead when a thread will be blocked for longer periods of time. The
6469 presence of both interfaces allows implementations with these two different characteristics,
6470 both of which may be useful to a particular application.

6471 It has also been suggested that applications can build their own spin locks from the
6472 *pthread_mutex_trylock()* function:

```
6473         while (pthread_mutex_trylock(&mutex));
```

6474 The apparent simplicity of this construct is somewhat deceiving, however. While the actual
6475 wait is quite efficient, various guarantees on the integrity of mutex objects (for example,
6476 priority inheritance rules) may add overhead to the successful path of the trylock operation
6477 that is not required of spin locks. One could, of course, add an attribute to the mutex to
6478 bypass such overhead, but the very act of finding and testing this attribute represents more
6479 overhead than is found in the typical spin lock.

6480 The need to hold spin lock overhead to an absolute minimum also makes it impossible to
6481 provide guarantees against starvation similar to those provided for mutexes or read-write
6482 locks. The overhead required to implement such guarantees (for example, disabling
6483 preemption before spinning) may well exceed the overhead of the spin wait itself by many
6484 orders of magnitude. If a “safe” spin wait seems desirable, it can always be provided (albeit
6485 at some performance cost) via appropriate mutex attributes.

6486 XSI Supported Functions

6487 On XSI-conformant systems, the following symbolic constants are always defined:

```
6488     _POSIX_READER_WRITER_LOCKS
6489     _POSIX_THREAD_ATTR_STACKADDR
6490     _POSIX_THREAD_ATTR_STACKSIZE
6491     _POSIX_THREAD_PROCESS_SHARED
6492     _POSIX_THREADS
```

6493 Therefore, the following threads functions are always supported:

6494	<i>pthread_atfork()</i>	<i>pthread_key_create()</i>
6495	<i>pthread_attr_destroy()</i>	<i>pthread_key_delete()</i>
6496	<i>pthread_attr_getdetachstate()</i>	<i>pthread_kill()</i>
6497	<i>pthread_attr_getguardsize()</i>	<i>pthread_mutex_destroy()</i>
6498	<i>pthread_attr_getschedparam()</i>	<i>pthread_mutex_init()</i>
6499	<i>pthread_attr_getstack()</i>	<i>pthread_mutex_lock()</i>
6500	<i>pthread_attr_getstackaddr()</i>	<i>pthread_mutex_trylock()</i>
6501	<i>pthread_attr_getstacksize()</i>	<i>pthread_mutex_unlock()</i>
6502	<i>pthread_attr_init()</i>	<i>pthread_mutexattr_destroy()</i>
6503	<i>pthread_attr_setdetachstate()</i>	<i>pthread_mutexattr_getpshared()</i>
6504	<i>pthread_attr_setguardsize()</i>	<i>pthread_mutexattr_gettype()</i>
6505	<i>pthread_attr_setschedparam()</i>	<i>pthread_mutexattr_init()</i>
6506	<i>pthread_attr_setstack()</i>	<i>pthread_mutexattr_setpshared()</i>
6507	<i>pthread_attr_setstackaddr()</i>	<i>pthread_mutexattr_settype()</i>
6508	<i>pthread_attr_setstacksize()</i>	<i>pthread_once()</i>
6509	<i>pthread_cancel()</i>	<i>pthread_rwlock_destroy()</i>
6510	<i>pthread_cleanup_pop()</i>	<i>pthread_rwlock_init()</i>
6511	<i>pthread_cleanup_push()</i>	<i>pthread_rwlock_rdlock()</i>
6512	<i>pthread_cond_broadcast()</i>	<i>pthread_rwlock_tryrdlock()</i>
6513	<i>pthread_cond_destroy()</i>	<i>pthread_rwlock_trywrlock()</i>
6514	<i>pthread_cond_init()</i>	<i>pthread_rwlock_unlock()</i>
6515	<i>pthread_cond_signal()</i>	<i>pthread_rwlock_wrlock()</i>
6516	<i>pthread_cond_timedwait()</i>	<i>pthread_rwlockattr_destroy()</i>
6517	<i>pthread_cond_wait()</i>	<i>pthread_rwlockattr_getpshared()</i>
6518	<i>pthread_condattr_destroy()</i>	<i>pthread_rwlockattr_init()</i>

6519	<i>pthread_condattr_getpshared()</i>	<i>pthread_rwlockattr_setpshared()</i>
6520	<i>pthread_condattr_init()</i>	<i>pthread_self()</i>
6521	<i>pthread_condattr_setpshared()</i>	<i>pthread_setcancelstate()</i>
6522	<i>pthread_create()</i>	<i>pthread_setcanceltype()</i>
6523	<i>pthread_detach()</i>	<i>pthread_setconcurrency()</i>
6524	<i>pthread_equal()</i>	<i>pthread_setspecific()</i>
6525	<i>pthread_exit()</i>	<i>pthread_sigmask()</i>
6526	<i>pthread_getconcurrency()</i>	<i>pthread_testcancel()</i>
6527	<i>pthread_getspecific()</i>	<i>sigwait()</i>
6528	<i>pthread_join()</i>	

6529 On XSI-conformant systems, the symbolic constant `_POSIX_THREAD_SAFE_FUNCTIONS` is
6530 always defined. Therefore, the following functions are always supported:

6531	<i>asctime_r()</i>	<i>getpwuid_r()</i>
6532	<i>ctime_r()</i>	<i>gmtime_r()</i>
6533	<i>flockfile()</i>	<i>localtime_r()</i>
6534	<i>ftrylockfile()</i>	<i>putc_unlocked()</i>
6535	<i>funlockfile()</i>	<i>putchar_unlocked()</i>
6536	<i>getc_unlocked()</i>	<i>rand_r()</i>
6537	<i>getchar_unlocked()</i>	<i>readdir_r()</i>
6538	<i>getgrgid_r()</i>	<i>strerror_r()</i>
6539	<i>getgrnam_r()</i>	<i>strtok_r()</i>
6540	<i>getpwnam_r()</i>	

6541 The following threads functions are only supported on XSI-conformant systems if the Realtime
6542 Threads Option Group is supported :

6543	<i>pthread_attr_getinheritsched()</i>	<i>pthread_mutex_getprioceiling()</i>
6544	<i>pthread_attr_getschedpolicy()</i>	<i>pthread_mutex_setprioceiling()</i>
6545	<i>pthread_attr_getscope()</i>	<i>pthread_mutexattr_getprioceiling()</i>
6546	<i>pthread_attr_setinheritsched()</i>	<i>pthread_mutexattr_getprotocol()</i>
6547	<i>pthread_attr_setschedpolicy()</i>	<i>pthread_mutexattr_setprioceiling()</i>
6548	<i>pthread_attr_setscope()</i>	<i>pthread_mutexattr_setprotocol()</i>
6549	<i>pthread_getschedparam()</i>	<i>pthread_setschedparam()</i>

6550 XSI Threads Extensions

6551 The following XSI extensions to POSIX.1c are now supported in IEEE Std 1003.1-200x as part of
6552 the alignment with the Single UNIX Specification:

- 6553 • Extended mutex attribute types
- 6554 • Read-write locks and attributes (also introduced by IEEE Std 1003.1j-2000 amendment)
- 6555 • Thread concurrency level
- 6556 • Thread stack guard size
- 6557 • Parallel I/O

6558 A total of 19 new functions were added.

6559 These extensions carefully follow the threads programming model specified in POSIX.1c. As
6560 with POSIX.1c, all the new functions return zero if successful; otherwise, an error number is

6561 returned to indicate the error.

6562 The concept of attribute objects was introduced in POSIX.1c to allow implementations to extend
6563 IEEE Std 1003.1-200x without changing the existing interfaces. Attribute objects were defined for
6564 threads, mutexes, and condition variables. Attributes objects are defined as implementation-
6565 defined opaque types to aid extensibility, and functions are defined to allow attributes to be set
6566 or retrieved. This model has been followed when adding the new type attribute of
6567 **pthread_mutexattr_t** or the new read-write lock attributes object **pthread_rwlockattr_t**.

6568 • Extended Mutex Attributes

6569 POSIX.1c defines a mutex attributes object as an implementation-defined opaque object of
6570 type **pthread_mutexattr_t**, and specifies a number of attributes which this object must have
6571 and a number of functions which manipulate these attributes. These attributes include
6572 *detachstate*, *inheritsched*, *schedparm*, *schedpolicy*, *contentionscope*, *stackaddr*, and *stacksize*.

6573 The System Interfaces volume of IEEE Std 1003.1-200x specifies another mutex attribute
6574 called *type*. The *type* attribute allows applications to specify the behavior of mutex locking
6575 operations in situations where the POSIX.1c behavior is undefined. The OSF DCE threads
6576 implementation, based on Draft 4 of POSIX.1c, specified a similar attribute. Note that the
6577 names of the attributes have changed somewhat from the OSF DCE threads implementation.

6578 The System Interfaces volume of IEEE Std 1003.1-200x also extends the specification of the
6579 following POSIX.1c functions which manipulate mutexes:

6580 *pthread_mutex_lock()*
6581 *pthread_mutex_trylock()*
6582 *pthread_mutex_unlock()*

6583 to take account of the new mutex attribute type and to specify behavior which was declared
6584 as undefined in POSIX.1c. How a calling thread acquires or releases a mutex now depends
6585 upon the mutex *type* attribute.

6586 The *type* attribute can have the following values:

6587 PTHREAD_MUTEX_NORMAL

6588 Basic mutex with no specific error checking built in. Does not report a deadlock error.

6589 PTHREAD_MUTEX_RECURSIVE

6590 Allows any thread to recursively lock a mutex. The mutex must be unlocked an equal
6591 number of times to release the mutex.

6592 PTHREAD_MUTEX_ERRORCHECK

6593 Detects and reports simple usage errors; that is, an attempt to unlock a mutex that is not
6594 locked by the calling thread or that is not locked at all, or an attempt to relock a mutex
6595 the thread already owns.

6596 PTHREAD_MUTEX_DEFAULT

6597 The default mutex type. May be mapped to any of the above mutex types or may be an
6598 implementation-defined type.

6599 *Normal* mutexes do not detect deadlock conditions; for example, a thread will hang if it tries
6600 to relock a normal mutex that it already owns. Attempting to unlock a mutex locked by
6601 another thread, or unlocking an unlocked mutex, results in undefined behavior. Normal
6602 mutexes will usually be the fastest type of mutex available on a platform but provide the
6603 least error checking.

6604 *Recursive* mutexes are useful for converting old code where it is difficult to establish clear
6605 boundaries of synchronization. A thread can relock a recursive mutex without first unlocking

6606 it. The relocking deadlock which can occur with normal mutexes cannot occur with this type
6607 of mutex. However, multiple locks of a recursive mutex require the same number of unlocks
6608 to release the mutex before another thread can acquire the mutex. Furthermore, this type of
6609 mutex maintains the concept of an owner. Thus, a thread attempting to unlock a recursive
6610 mutex which another thread has locked returns with an error. A thread attempting to unlock
6611 a recursive mutex that is not locked shall return with an error. Never use a recursive mutex
6612 with condition variables because the implicit unlock performed by *pthread_cond_wait()* or
6613 *pthread_cond_timedwait()* will not actually release the mutex if it had been locked multiple
6614 times.

6615 *Errorcheck* mutexes provide error checking and are useful primarily as a debugging aid. A
6616 thread attempting to relock an errorcheck mutex without first unlocking it returns with an
6617 error. Again, this type of mutex maintains the concept of an owner. Thus, a thread
6618 attempting to unlock an errorcheck mutex which another thread has locked returns with an
6619 error. A thread attempting to unlock an errorcheck mutex that is not locked also returns with
6620 an error. It should be noted that errorcheck mutexes will almost always be much slower than
6621 normal mutexes due to the extra state checks performed.

6622 The *default* mutex type provides implementation-defined error checking. The default mutex
6623 may be mapped to one of the other defined types or may be something entirely different.
6624 This enables each vendor to provide the mutex semantics which the vendor feels will be
6625 most useful to their target users. Most vendors will probably choose to make normal
6626 mutexes the default so as to give applications the benefit of the fastest type of mutexes
6627 available on their platform. Check your implementation's documentation.

6628 An application developer can use any of the mutex types almost interchangeably as long as
6629 the application does not depend upon the implementation detecting (or failing to detect) any
6630 particular errors. Note that a recursive mutex can be used with condition variable waits as
6631 long as the application never recursively locks the mutex.

6632 Two functions are provided for manipulating the *type* attribute of a mutex attributes object.
6633 This attribute is set or returned in the *type* parameter of these functions. The
6634 *pthread_mutexattr_settype()* function is used to set a specific type value while
6635 *pthread_mutexattr_gettype()* is used to return the type of the mutex. Setting the *type* attribute
6636 of a mutex attributes object affects only mutexes initialized using that mutex attributes
6637 object. Changing the *type* attribute does not affect mutexes previously initialized using that
6638 mutex attributes object.

6639 • Read-Write Locks and Attributes

6640 The read-write locks introduced have been harmonized with those in IEEE Std 1003.1j-2000;
6641 see also Section B.2.9.6 (on page 3464).

6642 Read-write locks (also known as reader-writer locks) allow a thread to exclusively lock some
6643 shared data while updating that data, or allow any number of threads to have simultaneous
6644 read-only access to the data.

6645 Unlike a mutex, a read-write lock distinguishes between reading data and writing data. A
6646 mutex excludes all other threads. A read-write lock allows other threads access to the data,
6647 providing no thread is modifying the data. Thus, a read-write lock is less primitive than
6648 either a mutex-condition variable pair or a semaphore.

6649 Application developers should consider using a read-write lock rather than a mutex to
6650 protect data that is frequently referenced but seldom modified. Most threads (readers) will be
6651 able to read the data without waiting and will only have to block when some other thread (a
6652 writer) is in the process of modifying the data. Conversely a thread that wants to change the
6653 data is forced to wait until there are no readers. This type of lock is often used to facilitate

6654 parallel access to data on multi-processor platforms or to avoid context switches on single
6655 processor platforms where multiple threads access the same data.

6656 If a read-write lock becomes unlocked and there are multiple threads waiting to acquire the
6657 write lock, the implementation's scheduling policy determines which thread shall acquire the
6658 read-write lock for writing. If there are multiple threads blocked on a read-write lock for both
6659 read locks and write locks, it is unspecified whether the readers or a writer acquire the lock
6660 first. However, for performance reasons, implementations often favor writers over readers to
6661 avoid potential writer starvation.

6662 A read-write lock object is an implementation-defined opaque object of type
6663 **pthread_rwlock_t** as defined in `<pthread.h>`. There are two different sorts of locks
6664 associated with a read-write lock: a *read lock* and a *write lock*.

6665 The `pthread_rwlockattr_init()` function initializes a read-write lock attributes object with the
6666 default value for all the attributes defined in the implementation. After a read-write lock
6667 attributes object has been used to initialize one or more read-write locks, changes to the
6668 read-write lock attributes object, including destruction, do not affect previously initialized
6669 read-write locks.

6670 Implementations must provide at least the read-write lock attribute *process-shared*. This
6671 attribute can have the following values:

6672 **PTHREAD_PROCESS_SHARED**
6673 Any thread of any process that has access to the memory where the read-write lock
6674 resides can manipulate the read-write lock.

6675 **PTHREAD_PROCESS_PRIVATE**
6676 Only threads created within the same process as the thread that initialized the read-
6677 write lock can manipulate the read-write lock. This is the default value.

6678 The `pthread_rwlockattr_setpshared()` function is used to set the *process-shared* attribute of an
6679 initialized read-write lock attributes object while the function `pthread_rwlockattr_getpshared()`
6680 obtains the current value of the *process-shared* attribute.

6681 A read-write lock attributes object is destroyed using the `pthread_rwlockattr_destroy()`
6682 function. The effect of subsequent use of the read-write lock attributes object is undefined.

6683 A thread creates a read-write lock using the `pthread_rwlock_init()` function. The attributes of
6684 the read-write lock can be specified by the application developer; otherwise, the default
6685 implementation-defined read-write lock attributes are used if the pointer to the read-write
6686 lock attributes object is NULL. In cases where the default attributes are appropriate, the
6687 **PTHREAD_RWLOCK_INITIALIZER** macro can be used to initialize statically allocated
6688 read-write locks.

6689 A thread which wants to apply a read lock to the read-write lock can use either
6690 `pthread_rwlock_rdlock()` or `pthread_rwlock_tryrdlock()`. If `pthread_rwlock_rdlock()` is used, the
6691 thread acquires a read lock if a writer does not hold the write lock and there are no writers
6692 blocked on the write lock. If a read lock is not acquired, the calling thread blocks until it can
6693 acquire a lock. However, if `pthread_rwlock_tryrdlock()` is used, the function returns
6694 immediately with the error [EBUSY] if any thread holds a write lock or there are blocked
6695 writers waiting for the write lock.

6696 A thread which wants to apply a write lock to the read-write lock can use either of two
6697 functions: `pthread_rwlock_wrlock()` or `pthread_rwlock_trywrlock()`. If `pthread_rwlock_wrlock()`
6698 is used, the thread acquires the write lock if no other reader or writer threads hold the read-
6699 write lock. If the write lock is not acquired, the thread blocks until it can acquire the write
6700 lock. However, if `pthread_rwlock_trywrlock()` is used, the function returns immediately with

6701 the error [EBUSY] if any thread is holding either a read or a write lock.

6702 The `pthread_rwlock_unlock()` function is used to unlock a read-write lock object held by the
6703 calling thread. Results are undefined if the read-write lock is not held by the calling thread. If
6704 there are other read locks currently held on the read-write lock object, the read-write lock
6705 object shall remain in the read locked state but without the current thread as one of its
6706 owners. If this function releases the last read lock for this read-write lock object, the read-
6707 write lock object shall be put in the unlocked read state. If this function is called to release a
6708 write lock for this read-write lock object, the read-write lock object shall be put in the
6709 unlocked state.

6710 • Thread Concurrency Level

6711 On threads implementations that multiplex user threads onto a smaller set of kernel
6712 execution entities, the system attempts to create a reasonable number of kernel execution
6713 entities for the application upon application startup.

6714 On some implementations, these kernel entities are retained by user threads that block in the
6715 kernel. Other implementations do not *timeslice* user threads so that multiple compute-bound
6716 user threads can share a kernel thread. On such implementations, some applications may use
6717 up all the available kernel execution entities before its user-space threads are used up. The
6718 process may be left with user threads capable of doing work for the application but with no
6719 way to schedule them.

6720 The `pthread_setconcurrency()` function enables an application to request more kernel entities;
6721 that is, specify a desired concurrency level. However, this function merely provides a hint to
6722 the implementation. The implementation is free to ignore this request or to provide some
6723 other number of kernel entities. If an implementation does not multiplex user threads onto a
6724 smaller number of kernel execution entities, the `pthread_setconcurrency()` function has no
6725 effect.

6726 The `pthread_setconcurrency()` function may also have an effect on implementations where the
6727 kernel mode and user mode schedulers cooperate to ensure that ready user threads are not
6728 prevented from running by other threads blocked in the kernel.

6729 The `pthread_getconcurrency()` function always returns the value set by a previous call to
6730 `pthread_setconcurrency()`. However, if `pthread_setconcurrency()` was not previously called, this
6731 function shall return zero to indicate that the threads implementation is maintaining the
6732 concurrency level.

6733 • Thread Stack Guard Size

6734 DCE threads introduced the concept of a *thread stack guard size*. Most thread
6735 implementations add a region of protected memory to a thread's stack, commonly known as
6736 a *guard region*, as a safety measure to prevent stack pointer overflow in one thread from
6737 corrupting the contents of another thread's stack. The default size of the guard regions
6738 attribute is {PAGESIZE} bytes and is implementation-defined.

6739 Some application developers may wish to change the stack guard size. When an application
6740 creates a large number of threads, the extra page allocated for each stack may strain system
6741 resources. In addition to the extra page of memory, the kernel's memory manager has to keep
6742 track of the different protections on adjoining pages. When this is a problem, the application
6743 developer may request a guard size of 0 bytes to conserve system resources by eliminating
6744 stack overflow protection.

6745 Conversely an application that allocates large data structures such as arrays on the stack may
6746 wish to increase the default guard size in order to detect stack overflow. If a thread allocates
6747 two pages for a data array, a single guard page provides little protection against thread stack

6748 overflows since the thread can corrupt adjoining memory beyond the guard page.

6749 The System Interfaces volume of IEEE Std 1003.1-200x defines a new attribute of a thread
6750 attributes object; that is, the *guardsize* attribute which allows applications to specify the size
6751 of the guard region of a thread's stack.

6752 Two functions are provided for manipulating a thread's stack guard size. The
6753 *pthread_attr_setguardsize()* function sets the thread *guardsize* attribute, and the
6754 *pthread_attr_getguardsize()* function retrieves the current value.

6755 An implementation may round up the requested guard size to a multiple of the configurable
6756 system variable {PAGESIZE}. In this case, *pthread_attr_getguardsize()* returns the guard size
6757 specified by the previous *pthread_attr_setguardsize()* function call and not the rounded up
6758 value.

6759 If an application is managing its own thread stacks using the *stackaddr* attribute, the *guardsize*
6760 attribute is ignored and no stack overflow protection is provided. In this case, it is the
6761 responsibility of the application to manage stack overflow along with stack allocation.

6762 • Parallel I/O

6763 Suppose two or more threads independently issue read requests on the same file. To read
6764 specific data from a file, a thread must first call *lseek()* to seek to the proper offset in the file,
6765 and then call *read()* to retrieve the required data. If more than one thread does this at the
6766 same time, the first thread may complete its seek call, but before it gets a chance to issue its
6767 read call a second thread may complete its seek call, resulting in the first thread accessing
6768 incorrect data when it issues its read call. One workaround is to lock the file descriptor while
6769 seeking and reading or writing, but this reduces parallelism and adds overhead.

6770 Instead, the System Interfaces volume of IEEE Std 1003.1-200x provides two functions to
6771 make seek/read and seek/write operations atomic. The file descriptor's current offset is
6772 unchanged, thus allowing multiple read and write operations to proceed in parallel. This
6773 improves the I/O performance of threaded applications. The *pread()* function is used to do
6774 an atomic read of data from a file into a buffer. Conversely, the *pwrite()* function does an
6775 atomic write of data from a buffer to a file.

6776 **B.2.9.1 Thread-Safety**

6777 All functions required by IEEE Std 1003.1-200x need to be thread-safe. Implementations have to
6778 provide internal synchronization when necessary in order to achieve this goal. In certain
6779 cases—for example, most floating-point implementations—context switch code may have to
6780 manage the writable shared state.

6781 While a read from a pipe of {PIPE_MAX}*2 bytes may not generate a single atomic and thread- |
6782 safe stream of bytes, it should generate “several” (individually atomic) thread-safe streams of |
6783 bytes. Similarly, while reading from a terminal device may not generate a single atomic and |
6784 thread-safe stream of bytes, it should generate some finite number of (individually atomic) and |
6785 thread-safe streams of bytes. That is, concurrent calls to read for a pipe, FIFO, or terminal device |
6786 are not allowed to result in corrupting the stream of bytes or other internal data. However, |
6787 *read()*, in these cases, is not required to return a single contiguous and atomic stream of bytes. |

6788 It is not required that all functions provided by IEEE Std 1003.1-200x be either async-cancel-safe |
6789 or async-signal-safe.

6790 As it turns out, some functions are inherently not thread-safe; that is, their interface
6791 specifications preclude reentrancy. For example, some functions (such as *asctime()*) return a
6792 pointer to a result stored in memory space allocated by the function on a per-process basis. Such
6793 a function is not thread-safe, because its result can be overwritten by successive invocations.

6794 Other functions, while not inherently non-thread-safe, may be implemented in ways that lead to
6795 them not being thread-safe. For example, some functions (such as *rand()*) store state information
6796 (such as a seed value, which survives multiple function invocations) in memory space allocated
6797 by the function on a per-process basis. The implementation of such a function is not thread-safe
6798 if the implementation fails to synchronize invocations of the function and thus fails to protect
6799 the state information. The problem is that when the state information is not protected,
6800 concurrent invocations can interfere with one another (for example, see the same seed value).

6801 *Thread-Safety and Locking of Existing Functions*

6802 Originally, POSIX.1 was not designed to work in a multi-threaded environment, and some
6803 implementations of some existing functions will not work properly when executed concurrently.
6804 To provide routines that will work correctly in an environment with threads (“thread-safe”), two
6805 problems need to be solved:

- 6806 1. Routines that maintain or return pointers to static areas internal to the routine (which may
6807 now be shared) need to be modified. The routines *ttyname()* and *localtime()* are examples.
- 6808 2. Routines that access data space shared by more than one thread need to be modified. The
6809 *malloc()* function and the *stdio* family routines are examples.

6810 There are a variety of constraints on these changes. The first is compatibility with the existing
6811 versions of these functions—non-thread-safe functions will continue to be in use for some time,
6812 as the original interfaces are used by existing code. Another is that the new thread-safe versions
6813 of these functions represent as small a change as possible over the familiar interfaces provided
6814 by the existing non-thread-safe versions. The new interfaces should be independent of any
6815 particular threads implementation. In particular, they should be thread-safe without depending
6816 on explicit thread-specific memory. Finally, there should be minimal performance penalty due to
6817 the changes made to the functions.

6818 It is intended that the list of functions from POSIX.1 that cannot be made thread-safe and for
6819 which corrected versions are provided be complete.

6820 *Thread-Safety and Locking Solutions*

6821 Many of the POSIX.1 functions were thread-safe and did not change at all. However, some
6822 functions (for example, the math functions typically found in **libm**) are not thread-safe because
6823 of writable shared global state. For instance, in IEEE Std 754-1985 floating-point
6824 implementations, the computation modes and flags are global and shared.

6825 Some functions are not thread-safe because a particular implementation is not reentrant,
6826 typically because of a non-essential use of static storage. These require only a new
6827 implementation.

6828 Thread-safe libraries are useful in a wide range of parallel (and asynchronous) programming
6829 environments, not just within pthreads. In order to be used outside the context of pthreads,
6830 however, such libraries still have to use some synchronization method. These could either be
6831 independent of the pthread synchronization operations, or they could be a subset of the pthread
6832 interfaces. Either method results in thread-safe library implementations that can be used without
6833 the rest of pthreads.

6834 Some functions, such as the *stdio* family interface and dynamic memory allocation functions
6835 such as *malloc()*, are interdependent routines that share resources (for example, buffers) across
6836 related calls. These require synchronization to work correctly, but they do not require any
6837 change to their external (user-visible) interfaces.

6838 In some cases, such as *getc()* and *putc()*, adding synchronization is likely to create an
6839 unacceptable performance impact. In this case, slower thread-safe synchronized functions are to

6840 be provided, but the original, faster (but unsafe) functions (which may be implemented as
6841 macros) are retained under new names. Some additional special-purpose synchronization
6842 facilities are necessary for these macros to be usable in multi-threaded programs. This also
6843 requires changes in `<stdio.h>`.

6844 The other common reason that functions are unsafe is that they return a pointer to static storage,
6845 making the functions non-thread-safe. This has to be changed, and there are three natural
6846 choices:

6847 1. Return a pointer to thread-specific storage

6848 This could incur a severe performance penalty on those architectures with a costly
6849 implementation of the thread-specific data interface.

6850 A variation on this technique is to use `malloc()` to allocate storage for the function output
6851 and return a pointer to this storage. This technique may also have an undesirable
6852 performance impact, however, and a simplistic implementation requires that the user
6853 program explicitly free the storage object when it is no longer needed. This technique is
6854 used by some existing POSIX.1 functions. With careful implementation for infrequently
6855 used functions, there may be little or no performance or storage penalty, and the
6856 maintenance of already-standardized interfaces is a significant benefit.

6857 2. Return the actual value computed by the function

6858 This technique can only be used with functions that return pointers to structures—routines
6859 that return character strings would have to wrap their output in an enclosing structure in
6860 order to return the output on the stack. There is also a negative performance impact
6861 inherent in this solution in that the output value has to be copied twice before it can be
6862 used by the calling function: once from the called routine's local buffers to the top of the
6863 stack, then from the top of the stack to the assignment target. Finally, many older
6864 compilers cannot support this technique due to a historical tendency to use internal static
6865 buffers to deliver the results of structure-valued functions.

6866 3. Have the caller pass the address of a buffer to contain the computed value

6867 The only disadvantage of this approach is that extra arguments have to be provided by the
6868 calling program. It represents the most efficient solution to the problem, however, and,
6869 unlike the `malloc()` technique, it is semantically clear.

6870 There are some routines (often groups of related routines) whose interfaces are inherently non-
6871 thread-safe because they communicate across multiple function invocations by means of static
6872 memory locations. The solution is to redesign the calls so that they are thread-safe, typically by
6873 passing the needed data as extra parameters. Unfortunately, this may require major changes to
6874 the interface as well.

6875 A floating-point implementation using IEEE Std 754-1985 is a case in point. A less problematic
6876 example is the `rand48` family of pseudo-random number generators. The functions `getgrgid()`,
6877 `getgrnam()`, `getpwnam()`, and `getpwuid()` are another such case.

6878 The problems with `errno` are discussed in **Alternative Solutions for Per-Thread `errno`** (on page
6879 3382).

6880 Some functions can be thread-safe or not, depending on their arguments. These include the
6881 `tmpnam()` and `ctermid()` functions. These functions have pointers to character strings as
6882 arguments. If the pointers are not NULL, the functions store their results in the character string;
6883 however, if the pointers are NULL, the functions store their results in an area that may be static
6884 and thus subject to overwriting by successive calls. These should only be called by multi-thread
6885 applications when their arguments are non-NULL.

6886 *Asynchronous Safety and Thread-Safety*

6887 A floating-point implementation has many modes that effect rounding and other aspects of
6888 computation. Functions in some math library implementations may change the computation
6889 modes for the duration of a function call. If such a function call is interrupted by a signal or
6890 cancelation, the floating-point state is not required to be protected.

6891 There is a significant cost to make floating-point operations async-cancel-safe or async-signal-
6892 safe; accordingly, neither form of async safety is required.

6893 *Functions Returning Pointers to Static Storage*

6894 For those functions that are not thread-safe because they return values in fixed size statically
6895 allocated structures, alternate “_r” forms are provided that pass a pointer to an explicit result
6896 structure. Those that return pointers into library-allocated buffers have forms provided with
6897 explicit buffer and length parameters.

6898 For functions that return pointers to library-allocated buffers, it makes sense to provide “_r”
6899 versions that allow the application control over allocation of the storage in which results are
6900 returned. This allows the state used by these functions to be managed on an application-specific
6901 basis, supporting per-thread, per-process, or other application-specific sharing relationships.

6902 Early proposals had provided “_r” versions for functions that returned pointers to variable-size
6903 buffers without providing a means for determining the required buffer size. This would have
6904 made using such functions exceedingly clumsy, potentially requiring iteratively calling them
6905 with increasingly larger guesses for the amount of storage required. Hence, *sysconf()* variables
6906 have been provided for such functions that return the maximum required buffer size.

6907 Thus, the rule that has been followed by IEEE Std 1003.1-200x when adapting single-threaded
6908 non-thread-safe functions is as follows: all functions returning pointers to library-allocated
6909 storage should have “_r” versions provided, allowing the application control over the storage
6910 allocation. Those with variable-sized return values accept both a buffer address and a length
6911 parameter. The *sysconf()* variables are provided to supply the appropriate buffer sizes when
6912 required. Implementors are encouraged to apply the same rule when adapting their own existing
6913 functions to a pthreads environment.

6914 *B.2.9.2 Thread IDs*

6915 Separate programs should communicate through well-defined interfaces and should not depend
6916 on each other's implementation. For example, if a programmer decides to rewrite the *sort*
6917 program using multiple threads, it should be easy to do this so that the interface to the *sort*
6918 program does not change. Consider that if the user causes SIGINT to be generated while the *sort*
6919 program is running, keeping the same interface means that the entire sort program is killed, not
6920 just one of its threads. As another example, consider a realtime program that manages a reactor.
6921 Such a program may wish to allow other programs to control the priority at which it watches the
6922 control rods. One technique to accomplish this is to write the ID of the thread watching the
6923 control rods into a file and allow other programs to change the priority of that thread as they see
6924 fit. A simpler technique is to have the reactor process accept IPCs (Inter-Process Communication
6925 messages) from other processes, telling it at a semantic level what priority the program should
6926 assign to watching the control rods. This allows the programmer greater flexibility in the
6927 implementation. For example, the programmer can change the implementation from having one
6928 thread per rod to having one thread watching all of the rods without changing the interface.
6929 Having threads live inside the process means that the implementation of a process is invisible to
6930 outside processes (excepting debuggers and system management tools).

6931 Threads do not provide a protection boundary. Every thread model allows threads to share
6932 memory with other threads and encourages this sharing to be widespread. This means that one

6933 thread can wipe out memory that is needed for the correct functioning of other threads that are
 6934 sharing its memory. Consequently, providing each thread with its own user and/or group IDs
 6935 would not provide a protection boundary between threads sharing memory.

6936 *B.2.9.3 Thread Mutexes*

6937 There is no additional rationale provided for this section.

6938 *B.2.9.4 Thread Scheduling*

6939 • Scheduling Implementation Models

6940 The following scheduling implementation models are presented in terms of threads and
 6941 “kernel entities”. This is to simplify exposition of the models, and it does not imply that an
 6942 implementation actually has an identifiable “kernel entity”.

6943 A kernel entity is not defined beyond the fact that it has scheduling attributes that are used to
 6944 resolve contention with other kernel entities for execution resources. A kernel entity may be
 6945 thought of as an envelope that holds a thread or a separate kernel thread. It is not a
 6946 conventional process, although it shares with the process the attribute that it has a single
 6947 thread of control; it does not necessarily imply an address space, open files, and so on. It is
 6948 better thought of as a primitive facility upon which conventional processes and threads may
 6949 be constructed.

6950 — System Thread Scheduling Model

6951 This model consists of one thread per kernel entity. The kernel entity is solely responsible
 6952 for scheduling thread execution on one or more processors. This model schedules all
 6953 threads against all other threads in the system using the scheduling attributes of the
 6954 thread.

6955 — Process Scheduling Model

6956 A generalized process scheduling model consists of two levels of scheduling. A threads
 6957 library creates a pool of kernel entities, as required, and schedules threads to run on them
 6958 using the scheduling attributes of the threads. Typically, the size of the pool is a function
 6959 of the simultaneously runnable threads, not the total number of threads. The kernel then
 6960 schedules the kernel entities onto processors according to their scheduling attributes,
 6961 which are managed by the threads library. This set model potentially allows a wide range
 6962 of mappings between threads and kernel entities.

6963 • System and Process Scheduling Model Performance

6964 There are a number of important implications on the performance of applications using these
 6965 scheduling models. The process scheduling model potentially provides lower overhead for
 6966 making scheduling decisions, since there is no need to access kernel-level information or
 6967 functions and the set of schedulable entities is smaller (only the threads within the process).

6968 On the other hand, since the kernel is also making scheduling decisions regarding the system
 6969 resources under its control (for example, CPU(s), I/O devices, memory), decisions that do
 6970 not take thread scheduling parameters into account can result in unspecified delays for
 6971 realtime application threads, causing them to miss maximum response time limits. |

6972 • Rate Monotonic Scheduling

6973 Rate monotonic scheduling was considered, but rejected for standardization in the context of
 6974 pthreads. A sporadic server policy is included.

6975 • Scheduling Options

6976 In IEEE Std 1003.1-200x, the basic thread scheduling functions are defined under the Threads
6977 option, so that they are required of all threads implementations. However, there are no
6978 specific scheduling policies required by this option to allow for conforming thread
6979 implementations that are not targeted to realtime applications.

6980 Specific standard scheduling policies are defined to be under the Thread Execution
6981 Scheduling option, and they are specifically designed to support realtime applications by
6982 providing predictable resource sharing sequences. The name of this option was chosen to
6983 emphasize that this functionality is defined as appropriate for realtime applications that
6984 require simple priority-based scheduling.

6985 It is recognized that these policies are not necessarily satisfactory for some multi-processor
6986 implementations, and work is ongoing to address a wider range of scheduling behaviors. The
6987 interfaces have been chosen to create abundant opportunity for future scheduling policies to
6988 be implemented and standardized based on this interface. In order to standardize a new
6989 scheduling policy, all that is required (from the standpoint of thread scheduling attributes) is
6990 to define a new policy name, new members of the thread attributes object, and functions to
6991 set these members when the scheduling policy is equal to the new value.

6992 **Scheduling Contention Scope**

6993 In order to accommodate the requirement for realtime response, each thread has a scheduling
6994 contention scope attribute. Threads with a system scheduling contention scope have to be
6995 scheduled with respect to all other threads in the system. These threads are usually bound to a
6996 single kernel entity that reflects their scheduling attributes and are directly scheduled by the
6997 kernel.

6998 Threads with a process scheduling contention scope need be scheduled only with respect to the
6999 other threads in the process. These threads may be scheduled within the process onto a pool of
7000 kernel entities. The implementation is also free to bind these threads directly to kernel entities
7001 and let them be scheduled by the kernel. Process scheduling contention scope allows the
7002 implementation the most flexibility and is the default if both contention scopes are supported
7003 and none is specified.

7004 Thus, the choice by implementors to provide one or the other (or both) of these scheduling
7005 models is driven by the need of their supported application domains for worst-case (that is,
7006 realtime) response, or average-case (non-realtime) response.

7007 **Scheduling Allocation Domain**

7008 The SCHED_FIFO and SCHED_RR scheduling policies take on different characteristics on a
7009 multi-processor. Other scheduling policies are also subject to changed behavior when executed
7010 on a multi-processor. The concept of scheduling allocation domain determines the set of
7011 processors on which the threads of an application may run. By considering the application's
7012 processor scheduling allocation domain for its threads, scheduling policies can be defined in
7013 terms of their behavior for varying processor scheduling allocation domain values. It is
7014 conceivable that not all scheduling allocation domain sizes make sense for all scheduling
7015 policies on all implementations. The concept of scheduling allocation domain, however, is a
7016 useful tool for the description of multi-processor scheduling policies.

7017 The “process control” approach to scheduling obtains significant performance advantages from
7018 dynamic scheduling allocation domain sizes when it is applicable.

7019 Non-Uniform Memory Access (NUMA) multi-processors may use a system scheduling structure
7020 that involves reassignment of threads among scheduling allocation domains. In NUMA

7021 machines, a natural model of scheduling is to match scheduling allocation domains to clusters of
7022 processors. Load balancing in such an environment requires changing the scheduling allocation
7023 domain to which a thread is assigned.

7024 **Scheduling Documentation**

7025 Implementation-provided scheduling policies need to be completely documented in order to be
7026 useful. This documentation includes a description of the attributes required for the policy, the
7027 scheduling interaction of threads running under this policy and all other supported policies, and
7028 the effects of all possible values for processor scheduling allocation domain. Note that for the
7029 implementor wishing to be minimally-compliant, it is (minimally) acceptable to define the
7030 behavior as undefined.

7031 **Scheduling Contention Scope Attribute**

7032 The scheduling contention scope defines how threads compete for resources. Within
7033 IEEE Std 1003.1-200x, scheduling contention scope is used to describe only how threads are
7034 scheduled in relation to one another in the system. That is, either they are scheduled against all
7035 other threads in the system (“system scope”) or only against those threads in the process
7036 (“process scope”). In fact, scheduling contention scope may apply to additional resources,
7037 including virtual timers and profiling, which are not currently considered by
7038 IEEE Std 1003.1-200x.

7039 **Mixed Scopes**

7040 If only one scheduling contention scope is supported, the scheduling decision is straightforward.
7041 To perform the processor scheduling decision in a mixed scope environment, it is necessary to
7042 map the scheduling attributes of the thread with process-wide contention scope to the same
7043 attribute space as the thread with system-wide contention scope.

7044 Since a conforming implementation has to support one and may support both scopes, it is useful
7045 to discuss the effects of such choices with respect to example applications. If an implementation
7046 supports both scopes, mixing scopes provides a means of better managing system-level (that is,
7047 kernel-level) and library-level resources. In general, threads with system scope will require the
7048 resources of a separate kernel entity in order to guarantee the scheduling semantics. On the
7049 other hand, threads with process scope can share the resources of a kernel entity while
7050 maintaining the scheduling semantics.

7051 The application is free to create threads with dedicated kernel resources, and other threads that
7052 multiplex kernel resources. Consider the example of a window server. The server allocates two
7053 threads per widget: one thread manages the widget user interface (including drawing), while the
7054 other thread takes any required application action. This allows the widget to be “active” while
7055 the application is computing. A screen image may be built from thousands of widgets. If each of
7056 these threads had been created with system scope, then most of the kernel-level resources might
7057 be wasted, since only a few widgets are active at any one time. In addition, mixed scope is
7058 particularly useful in a window server where one thread with high priority and system scope
7059 handles the mouse so that it tracks well. As another example, consider a database server. For
7060 each of the hundreds or thousands of clients supported by a large server, an equivalent number
7061 of threads will have to be created. If each of these threads were system, the consequences would
7062 be the same as for the window server example above. However, the server could be constructed
7063 so that actual retrieval of data is done by several dedicated threads. Dedicated threads that do
7064 work for all clients frequently justify the added expense of system scope. If it were not
7065 permissible to mix system and process threads in the same process, this type of solution would
7066 not be possible.

7067 Dynamic Thread Scheduling Parameters Access

7068 In many time-constrained applications, there is no need to change the scheduling attributes
7069 dynamically during thread or process execution, since the general use of these attributes is to
7070 reflect directly the time constraints of the application. Since these time constraints are generally
7071 imposed to meet higher-level system requirements, such as accuracy or availability, they
7072 frequently should remain unchanged during application execution.

7073 However, there are important situations in which the scheduling attributes should be changed.
7074 Generally, this will occur when external environmental conditions exist in which the time
7075 constraints change. Consider, for example, a space vehicle major mode change, such as the
7076 change from ascent to descent mode, or the change from the space environment to the
7077 atmospheric environment. In such cases, the frequency with which many of the sensors or
7078 acutators need to be read or written will change, which will necessitate a priority change. In
7079 other cases, even the existence of a time constraint might be temporary, necessitating not just a
7080 priority change, but also a policy change for ongoing threads or processes. For this reason, it is
7081 critical that the interface should provide functions to change the scheduling parameters
7082 dynamically, but, as with many of the other realtime functions, it is important that applications
7083 use them properly to avoid the possibility of unnecessarily degrading performance.

7084 In providing functions for dynamically changing the scheduling behavior of threads, there were
7085 two options: provide functions to get and set the individual scheduling parameters of threads, or
7086 provide a single interface to get and set all the scheduling parameters for a given thread
7087 simultaneously. Both approaches have merit. Access functions for individual parameters allow
7088 simpler control of thread scheduling for simple thread scheduling parameters. However, a single
7089 function for setting all the parameters for a given scheduling policy is required when first setting
7090 that scheduling policy. Since the single all-encompassing functions are required, it was decided
7091 to leave the interface as minimal as possible. Note that simpler functions (such as
7092 *pthread_setprio()* for threads running under the priority-based schedulers) can be easily defined
7093 in terms of the all-encompassing functions.

7094 If the *pthread_setschedparam()* function executes successfully, it will have set all of the scheduling
7095 parameter values indicated in *param*; otherwise, none of the scheduling parameters will have
7096 been modified. This is necessary to ensure that the scheduling of this and all other threads
7097 continues to be consistent in the presence of an erroneous scheduling parameter.

7098 The [EPERM] error value is included in the list of possible *pthread_setschedparam()* error returns
7099 as a reflection of the fact that the ability to change scheduling parameters increases risks to the
7100 implementation and application performance if the scheduling parameters are changed
7101 improperly. For this reason, and based on some existing practice, it was felt that some
7102 implementations would probably choose to define specific permissions for changing either a
7103 thread's own or another thread's scheduling parameters. IEEE Std 1003.1-200x does not include
7104 portable methods for setting or retrieving permissions, so any such use of permissions is
7105 completely unspecified .

7106 Mutex Initialization Scheduling Attributes

7107 In a priority-driven environment, a direct use of traditional primitives like mutexes and
7108 condition variables can lead to unbounded priority inversion, where a higher priority thread can
7109 be blocked by a lower priority thread, or set of threads, for an unbounded duration of time. As a
7110 result, it becomes impossible to guarantee thread deadlines. Priority inversion can be bounded
7111 and minimized by the use of priority inheritance protocols. This allows thread deadlines to be
7112 guaranteed even in the presence of synchronization requirements.

7113 Two useful but simple members of the family of priority inheritance protocols are the basic
7114 priority inheritance protocol and the priority ceiling protocol emulation. Under the Basic Priority

7115 Inheritance protocol (governed by the Thread Priority Inheritance option), a thread that is
7116 blocking higher priority threads executes at the priority of the highest priority thread that it
7117 blocks. This simple mechanism allows priority inversion to be bounded by the duration of
7118 critical sections and makes timing analysis possible.

7119 Under the Priority Ceiling Protocol Emulation protocol (governed by the Thread Priority
7120 Protection option), each mutex has a priority ceiling, usually defined as the priority of the
7121 highest priority thread that can lock the mutex. When a thread is executing inside critical
7122 sections, its priority is unconditionally increased to the highest of the priority ceilings of all the
7123 mutexes owned by the thread. This protocol has two very desirable properties in uni-processor
7124 systems. First, a thread can be blocked by a lower priority thread for at most the duration of one
7125 single critical section. Furthermore, when the protocol is correctly used in a single processor, and
7126 if threads do not become blocked while owning mutexes, mutual deadlocks are prevented.

7127 The priority ceiling emulation can be extended to multiple processor environments, in which
7128 case the values of the priority ceilings will be assigned depending on the kind of mutex that is
7129 being used: local to only one processor, or global, shared by several processors. Local priority
7130 ceilings will be assigned the usual way, equal to the priority of the highest priority thread that
7131 may lock that mutex. Global priority ceilings will usually be assigned a priority level higher than
7132 all the priorities assigned to any of the threads that reside in the involved processors to avoid the
7133 effect called remote blocking.

7134 **Change the Priority Ceiling of a Mutex**

7135 In order for the priority protect protocol to exhibit its desired properties of bounding priority
7136 inversion and avoidance of deadlock, it is critical that the ceiling priority of a mutex be the same
7137 as the priority of the highest thread that can ever hold it, or higher. Thus, if the priorities of the
7138 threads using such mutexes never change dynamically, there is no need ever to change the
7139 priority ceiling of a mutex.

7140 However, if a major system mode change results in an altered response time requirement for one
7141 or more application threads, their priority has to change to reflect it. It will occasionally be the
7142 case that the priority ceilings of mutexes held also need to change. While changing priority
7143 ceilings should generally be avoided, it is important that IEEE Std 1003.1-200x provide these
7144 interfaces for those cases in which it is necessary.

7145 **B.2.9.5 Thread Cancellation**

7146 Many existing threads packages have facilities for canceling an operation or canceling a thread.
7147 These facilities are used for implementing user requests (such as the CANCEL button in a
7148 window-based application), for implementing OR parallelism (for example, telling the other
7149 threads to stop working once one thread has found a forced mate in a parallel chess program), or
7150 for implementing the ABORT mechanism in Ada.

7151 POSIX programs traditionally have used the signal mechanism combined with either *longjmp()*
7152 or polling to cancel operations. Many POSIX programmers have trouble using these facilities to
7153 solve their problems efficiently in a single-threaded process. With the introduction of threads,
7154 these solutions become even more difficult to use.

7155 The main issues with implementing a cancellation facility are specifying the operation to be
7156 canceled, cleanly releasing any resources allocated to that operation, controlling when the target
7157 notices that it has been canceled, and defining the interaction between asynchronous signals and
7158 cancellation.

7159 Specifying the Operation to Cancel

7160 Consider a thread that calls through five distinct levels of program abstraction and then, inside
7161 the lowest-level abstraction, calls a function that suspends the thread. (An abstraction boundary
7162 is a layer at which the client of the abstraction sees only the service being provided and can
7163 remain ignorant of the implementation. Abstractions are often layered, each level of abstraction
7164 being a client of the lower-level abstraction and implementing a higher-level abstraction.)
7165 Depending on the semantics of each abstraction, one could imagine wanting to cancel only the
7166 call that causes suspension, only the bottom two levels, or the operation being done by the entire
7167 thread. Canceling operations at a finer grain than the entire thread is difficult because threads
7168 are active and they may be run in parallel on a multi-processor. By the time one thread can make
7169 a request to cancel an operation, the thread performing the operation may have completed that
7170 operation and gone on to start another operation whose cancelation is not desired. Thread IDs
7171 are not reused until the thread has exited, and either it was created with the *Attr detachstate*
7172 attribute set to *PTHREAD_CREATE_DETACHED* or the *pthread_join()* or *pthread_detach()*
7173 function has been called for that thread. Consequently, a thread cancelation will never be
7174 misdirected when the thread terminates. For these reasons, the canceling of operations is done at
7175 the granularity of the thread. Threads are designed to be inexpensive enough so that a separate
7176 thread may be created to perform each separately cancelable operation; for example, each
7177 possibly long running user request.

7178 For cancelation to be used in existing code, cancelation scopes and handlers will have to be
7179 established for code that needs to release resources upon cancelation, so that it follows the
7180 programming discipline described in the text.

7181 A Special Signal Versus a Special Interface

7182 Two different mechanisms were considered for providing the cancelation interfaces. The first
7183 was to provide an interface to direct signals at a thread and then to define a special signal that
7184 had the required semantics. The other alternative was to use a special interface that delivered the
7185 correct semantics to the target thread.

7186 The solution using signals produced a number of problems. It required the implementation to
7187 provide cancelation in terms of signals whereas a perfectly valid (and possibly more efficient)
7188 implementation could have both layered on a low-level set of primitives. There were so many
7189 exceptions to the special signal (it cannot be used with kill, no POSIX.1 interfaces can be used
7190 with it) that it was clearly not a valid signal. Its semantics on delivery were also completely
7191 different from any existing POSIX.1 signal. As such, a special interface that did not mandate the
7192 implementation and did not confuse the semantics of signals and cancelation was felt to be the
7193 better solution.

7194 Races Between Cancelation and Resuming Execution

7195 Due to the nature of cancelation, there is generally no synchronization between the thread
7196 requesting the cancelation of a blocked thread and events that may cause that thread to resume
7197 execution. For this reason, and because excess serialization hurts performance, when both an
7198 event that a thread is waiting for has occurred and a cancelation request has been made and
7199 cancelation is enabled, IEEE Std 1003.1-200x explicitly allows the implementation to choose
7200 between returning from the blocking call or acting on the cancelation request.

7201 **Interaction of Cancellation with Asynchronous Signals**

7202 A typical use of cancellation is to acquire a lock on some resource and to establish a cancellation
7203 cleanup handler for releasing the resource when and if the thread is canceled.

7204 A correct and complete implementation of cancellation in the presence of asynchronous signals
7205 requires considerable care. An implementation has to push a cancellation cleanup handler on the
7206 cancellation cleanup stack while maintaining the integrity of the stack data structure. If an
7207 asynchronously generated signal is posted to the thread during a stack operation, the signal
7208 handler cannot manipulate the cancellation cleanup stack. As a consequence, asynchronous
7209 signal handlers may not cancel threads or otherwise manipulate the cancellation state of a thread.
7210 Threads may, of course, be canceled by another thread that used a *sigwait()* function to wait
7211 synchronously for an asynchronous signal.

7212 In order for cancellation to function correctly, it is required that asynchronous signal handlers not
7213 change the cancellation state. This requires that some elements of existing practice, such as using
7214 *longjmp()* to exit from an asynchronous signal handler implicitly, be prohibited in cases where
7215 the integrity of the cancellation state of the interrupt thread cannot be ensured.

7216 **Thread Cancellation Overview**

7217 • Cancellability States

7218 The three possible cancellability states (disabled, deferred, and asynchronous) are encoded
7219 into two separate bits ((disable, enable) and (deferred, asynchronous)) to allow them to be
7220 changed and restored independently. For instance, short code sequences that will not block
7221 sometimes disable cancellability on entry and restore the previous state upon exit. Likewise,
7222 long or unbounded code sequences containing no convenient explicit cancellation points will
7223 sometimes set the cancellability type to asynchronous on entry and restore the previous value
7224 upon exit.

7225 • Cancellation Points

7226 Cancellation points are points inside of certain functions where a thread has to act on any
7227 pending cancellation request when cancellability is enabled, if the function would block. As
7228 with checking for signals, operations need only check for pending cancellation requests when
7229 the operation is about to block indefinitely.

7230 The idea was considered of allowing implementations to define whether blocking calls such
7231 as *read()* should be cancellation points. It was decided that it would adversely affect the
7232 design of conforming applications if blocking calls were not cancellation points because
7233 threads could be left blocked in an uncancelable state.

7234 There are several important blocking routines that are specifically not made cancellation
7235 points:

7236 — *pthread_mutex_lock()*

7237 If *pthread_mutex_lock()* were a cancellation point, every routine that called it would also
7238 become a cancellation point (that is, any routine that touched shared state would
7239 automatically become a cancellation point). For example, *malloc()*, *free()*, and *rand()*
7240 would become cancellation points under this scheme. Having too many cancellation points
7241 makes programming very difficult, leading to either much disabling and restoring of
7242 cancellability or much difficulty in trying to arrange for reliable cleanup at every possible
7243 place.

7244 Since *pthread_mutex_lock()* is not a cancellation point, threads could result in being
7245 blocked uninterruptibly for long periods of time if mutexes were used as a general

7246 synchronization mechanism. As this is normally not acceptable, mutexes should only be
7247 used to protect resources that are held for small fixed lengths of time where not being
7248 able to be canceled will not be a problem. Resources that need to be held exclusively for
7249 long periods of time should be protected with condition variables.

7250 — *barrier_wait()*

7251 Canceling a barrier wait will render a barrier unusable. Similar to a barrier timeout (which
7252 the standard developers rejected), there is no way to guarantee the consistency of a
7253 barrier's internal data structures if a barrier wait is canceled.

7254 — *pthread_spin_lock()*

7255 As with mutexes, spin locks should only be used to protect resources that are held for
7256 small fixed lengths of time where not being cancelable will not be a problem.

7257 Every library routine should specify whether or not it includes any cancellation points.
7258 Typically, only those routines that may block or compute indefinitely need to include
7259 cancellation points.

7260 Correctly coded routines only reach cancellation points after having set up a cancellation
7261 cleanup handler to restore invariants if the thread is canceled at that point. Being cancelable
7262 only at specified cancellation points allows programmers to keep track of actions needed in a
7263 cancellation cleanup handler more easily. A thread should only be made asynchronously
7264 cancelable when it is not in the process of acquiring or releasing resources or otherwise in a
7265 state from which it would be difficult or impossible to recover.

7266 • Thread Cancellation Cleanup Handlers

7267 The cancellation cleanup handlers provide a portable mechanism, easy to implement, for
7268 releasing resources and restoring invariants. They are easier to use than signal handlers
7269 because they provide a stack of cancellation cleanup handlers rather than a single handler,
7270 and because they have an argument that can be used to pass context information to the
7271 handler.

7272 The alternative to providing these simple cancellation cleanup handlers (whose only use is for
7273 cleaning up when a thread is canceled) is to define a general exception package that could be
7274 used for handling and cleaning up after hardware traps and software detected errors. This
7275 was too far removed from the charter of providing threads to handle asynchrony. However,
7276 it is an explicit goal of IEEE Std 1003.1-200x to be compatible with existing exception facilities
7277 and languages having exceptions.

7278 The interaction of this facility and other procedure-based or language-level exception
7279 facilities is unspecified in this version of IEEE Std 1003.1-200x. However, it is intended that it
7280 be possible for an implementation to define the relationship between these cancellation
7281 cleanup handlers and Ada, C++, or other language-level exception handling facilities.

7282 It was suggested that the cancellation cleanup handlers should also be called when the
7283 process exits or calls the *exec* function. This was rejected partly due to the performance
7284 problem caused by having to call the cancellation cleanup handlers of every thread before the
7285 operation could continue. The other reason was that the only state expected to be cleaned up
7286 by the cancellation cleanup handlers would be the intraprocess state. Any handlers that are to
7287 clean up the interprocess state would be registered with *atexit()*. There is the orthogonal
7288 problem that the *exec* functions do not honor the *atexit()* handlers, but resolving this is
7289 beyond the scope of IEEE Std 1003.1-200x.

7290 • Async-Cancel Safety

7291 A function is said to be *async-cancel safe* if it is written in such a way that entering the function
7292 with asynchronous cancelability enabled will not cause any invariants to be violated, even if
7293 a cancelation request is delivered at any arbitrary instruction. Functions that are *async-*
7294 *cancel-safe* are often written in such a way that they need to acquire no resources for their
7295 operation and the visible variables that they may write are strictly limited.

7296 Any routine that gets a resource as a side-effect cannot be made *async-cancel-safe* (for
7297 example, *malloc()*). If such a routine were called with asynchronous cancelability enabled, it
7298 might acquire the resource successfully, but as it was returning to the client, it could act on a
7299 cancelation request. In such a case, the application would have no way of knowing whether
7300 the resource was acquired or not.

7301 Indeed, because many interesting routines cannot be made *async-cancel-safe*, most library
7302 routines in general are not *async-cancel-safe*. Every library routine should specify whether or
7303 not it is *async-cancel safe* so that programmers know which routines can be called from code
7304 that is asynchronously cancelable.

7305 **B.2.9.6 Thread Read-Write Locks**7306 **Background**

7307 Read-write locks are often used to allow parallel access to data on multi-processors, to avoid
7308 context switches on uni-processors when multiple threads access the same data, and to protect
7309 data structures that are frequently accessed (that is, read) but rarely updated (that is, written).
7310 The in-core representation of a file system directory is a good example of such a data structure.
7311 One would like to achieve as much concurrency as possible when searching directories, but limit
7312 concurrent access when adding or deleting files.

7313 Although read-write locks can be implemented with mutexes and condition variables, such
7314 implementations are significantly less efficient than is possible. Therefore, this synchronization
7315 primitive is included in IEEE Std 1003.1-200x for the purpose of allowing more efficient
7316 implementations in multi-processor systems.

7317 **Queuing of Waiting Threads**

7318 The *pthread_rwlock_unlock()* function description states that one writer or one or more readers
7319 shall acquire the lock if it is no longer held by any thread as a result of the call. However, the
7320 function does not specify which thread(s) acquire the lock, unless the Thread Execution
7321 Scheduling option is supported.

7322 The standard developers considered the issue of scheduling with respect to the queuing of
7323 threads blocked on a read-write lock. The question turned out to be whether
7324 IEEE Std 1003.1-200x should require priority scheduling of read-write locks for threads whose
7325 execution scheduling policy is priority-based (for example, *SCHED_FIFO* or *SCHED_RR*). There
7326 are tradeoffs between priority scheduling, the amount of concurrency achievable among readers,
7327 and the prevention of writer and/or reader starvation.

7328 For example, suppose one or more readers hold a read-write lock and the following threads
7329 request the lock in the listed order:

7330 pthread_rwlock_wrlock() - Low priority thread writer_a
 7331 pthread_rwlock_rdlock() - High priority thread reader_a
 7332 pthread_rwlock_rdlock() - High priority thread reader_b
 7333 pthread_rwlock_rdlock() - High priority thread reader_c

7334 When the lock becomes available, should *writer_a* block the high priority readers? Or, suppose a
 7335 read-write lock becomes available and the following are queued:

7336 pthread_rwlock_rdlock() - Low priority thread reader_a
 7337 pthread_rwlock_rdlock() - Low priority thread reader_b
 7338 pthread_rwlock_rdlock() - Low priority thread reader_c
 7339 pthread_rwlock_wrlock() - Medium priority thread writer_a
 7340 pthread_rwlock_rdlock() - High priority thread reader_d

7341 If priority scheduling is applied then *reader_d* would acquire the lock and *writer_a* would block
 7342 the remaining readers. But should the remaining readers also acquire the lock to increase
 7343 concurrency? The solution adopted takes into account that when the Thread Execution
 7344 Scheduling option is supported, high priority threads may in fact starve low priority threads (the
 7345 application developer is responsible in this case to design the system in such a way that this
 7346 starvation is avoided). Therefore, IEEE Std 1003.1-200x specifies that high priority readers take
 7347 precedence over lower priority writers. However, to prevent writer starvation from threads of
 7348 the same or lower priority, writers take precedence over readers of the same or lower priority.

7349 Priority inheritance mechanisms are non-trivial in the context of read-write locks. When a high
 7350 priority writer is forced to wait for multiple readers, for example, it is not clear which subset of
 7351 the readers should inherit the writer's priority. Furthermore, the internal data structures that
 7352 record the inheritance must be accessible to all readers, and this implies some sort of
 7353 serialization that could negate any gain in parallelism achieved through the use of multiple
 7354 readers in the first place. Finally, existing practice does not support the use of priority
 7355 inheritance for read-write locks. Therefore, no specification of priority inheritance or priority
 7356 ceiling is attempted. If reliable priority-scheduled synchronization is absolutely required, it can
 7357 always be obtained through the use of mutexes.

7358 **Comparison to *fcntl()* Locks**

7359 The read-write locks and the *fcntl()* locks in IEEE Std 1003.1-200x share a common goal:
 7360 increasing concurrency among readers, thus increasing throughput and decreasing delay.

7361 However, the read-write locks have two features not present in the *fcntl()* locks. First, under
 7362 priority scheduling, read-write locks are granted in priority order. Second, also under priority
 7363 scheduling, writer starvation is prevented by giving writers preference over readers of equal or
 7364 lower priority.

7365 Also, read-write locks can be used in systems lacking a file system, such as those conforming to
 7366 the minimal realtime system profile of IEEE Std 1003.13-1998.

7367 **History of Resolution Issues**

7368 Based upon some balloting objections, the draft specified the behavior of threads waiting on a
 7369 read-write lock during the execution of a signal handler, as if the thread had not called the lock
 7370 operation. However, this specified behavior would require implementations to establish
 7371 internal signal handlers even though this situation would be rare, or never happen for many
 7372 programs. This would introduce an unacceptable performance hit in comparison to the little
 7373 additional functionality gained. Therefore, the behavior of read-write locks and signals was
 7374 reverted back to its previous mutex-like specification.

7375 **B.2.9.7** *Thread Interactions with Regular File Operations*

7376 There is no additional rationale provided for this section.

7377 **B.2.10** **Sockets**

7378 The base document for the sockets interfaces in IEEE Std 1003.1-200x is the XNS, Issue 5.2
7379 specification. This was primarily chosen as it aligns with IPv6. Additional material has been
7380 added from IEEE Std 1003.1g-2000, notably socket concepts, raw sockets, the *pselect()* function,
7381 and the `<sys/select.h>` header.

7382 **B.2.10.1** *Address Families*

7383 There is no additional rationale provided for this section.

7384 **B.2.10.2** *Addressing*

7385 There is no additional rationale provided for this section.

7386 **B.2.10.3** *Protocols*

7387 There is no additional rationale provided for this section.

7388 **B.2.10.4** *Routing*

7389 There is no additional rationale provided for this section.

7390 **B.2.10.5** *Interfaces*

7391 There is no additional rationale provided for this section.

7392 **B.2.10.6** *Socket Types*

7393 The type `socklen_t` was invented to cover the range of implementations seen in the field. The
7394 intent of `socklen_t` is to be the type for all lengths that are naturally bounded in size; that is, that
7395 they are the length of a buffer which cannot sensibly become of massive size: network addresses,
7396 host names, string representations of these, ancillary data, control messages, and socket options
7397 are examples. Truly boundless sizes are represented by `size_t` as in *read()*, *write()*, and so on.

7398 All `socklen_t` types were originally (in BSD UNIX) of type `int`. During the development of
7399 IEEE Std 1003.1-200x, it was decided to change all buffer lengths to `size_t`, which appears at face
7400 value to make sense. When dual mode 32/64-bit systems came along, this choice unnecessarily
7401 complicated system interfaces because `size_t` (with `long`) was a different size under ILP32 and
7402 LP64 models. Reverting to `int` would have happened except that some implementations had
7403 already shipped 64-bit-only interfaces. The compromise was a type which could be defined to be
7404 any size by the implementation: `socklen_t`.

7405 **B.2.10.7** *Socket I/O Mode*

7406 There is no additional rationale provided for this section.

7407 *B.2.10.8 Socket Owner*

7408 There is no additional rationale provided for this section.

7409 *B.2.10.9 Socket Queue Limits*

7410 There is no additional rationale provided for this section.

7411 *B.2.10.10 Pending Error*

7412 There is no additional rationale provided for this section.

7413 *B.2.10.11 Socket Receive Queue*

7414 There is no additional rationale provided for this section.

7415 *B.2.10.12 Socket Out-of-Band Data State*

7416 There is no additional rationale provided for this section.

7417 *B.2.10.13 Connection Indication Queue*

7418 There is no additional rationale provided for this section.

7419 *B.2.10.14 Signals*

7420 There is no additional rationale provided for this section.

7421 *B.2.10.15 Asynchronous Errors*

7422 There is no additional rationale provided for this section.

7423 *B.2.10.16 Use of Options*

7424 There is no additional rationale provided for this section.

7425 *B.2.10.17 Use of Sockets for Local UNIX Connections*

7426 There is no additional rationale provided for this section.

7427 *B.2.10.18 Use of Sockets over Internet Protocols*

7428 A raw socket allows privileged users direct access to a protocol; for example, raw access to the
7429 IP and ICMP protocols is possible through raw sockets. Raw sockets are intended for
7430 knowledgeable applications that wish to take advantage of some protocol feature not directly
7431 accessible through the other sockets interfaces.

7432 *B.2.10.19 Use of Sockets over Internet Protocols Based on IPv4*

7433 There is no additional rationale provided for this section.

7434 *B.2.10.20 Use of Sockets over Internet Protocols Based on IPv6*

7435 There is no additional rationale provided for this section.

7436 **B.2.11 Tracing**

7437 The organization of the tracing rationale differs from the traditional rationale in that this tracing
7438 rationale text is written against the trace interface as a whole, rather than against the individual
7439 components of the trace interface or the normative section in which those components are
7440 defined. Therefore the sections below do not parallel the sections of normative text in
7441 IEEE Std 1003.1-200x.

7442 *B.2.11.1 Objectives*

7443 The intended uses of tracing are application-system debugging during system development, as a
7444 “flight recorder” for maintenance of fielded systems, and as a performance measurement tool. In
7445 all of these intended uses, the vendor-supplied computer system and its software are, for this
7446 discussion, assumed error-free; the intent being to debug the user-written and/or third-party
7447 application code, and their interactions. Clearly, problems with the vendor-supplied system and
7448 its software will be uncovered from time to time, but this is a byproduct of the primary activity,
7449 debugging user code.

7450 Another need for defining a trace interface in POSIX stems from the objective to provide an
7451 efficient portable way to perform benchmarks. Existing practice shows that such interfaces are
7452 commonly used in a variety of systems but with little commonality. As part of the benchmarking
7453 needs, we must consider two aspects within the trace interface.

7454 The first, and perhaps more important one, is the qualitative aspect.

7455 The second is the quantitative aspect.

7456 • Qualitative Aspect

7457 To better understand this aspect, let us consider an example. Suppose that you want to
7458 organize a number of actions to be performed during the day. Some of these actions are
7459 known at the beginning of the day. Some others, which may be more or less important, will
7460 be triggered by reading your mail. During the day you will make some phone calls and
7461 synchronously receive some more information. Finally you will receive asynchronous phone
7462 calls that also will trigger actions. If you, or somebody else, examines your day at work, you,
7463 or he, can discover that you have not efficiently organized your work. For instance, relative
7464 to the phone calls you made, would it be preferable to make some of these early in the
7465 morning? Or to delay some others until the end of the day? Relative to the phone calls you
7466 have received, you might find that somebody you called in the morning has called you 10
7467 times while you were performing some important work. To examine, afterwards, your day at
7468 work, you record in sequence all the trace events relative to your work. This should give you
7469 a chance of organizing your next day at work.

7470 This is the qualitative aspect of the trace interface. The user of a system needs to keep a trace
7471 of particular points the application passes through, so that he can eventually make some
7472 changes in the application and/or system configuration, to give the application a chance of
7473 running more efficiently.

7474 • Quantitative Aspect

7475 This aspect concerns primarily realtime applications, where missed deadlines can be
7476 undesirable. Although there are, in IEEE Std 1003.1-200x, some interfaces useful for such
7477 applications (timeouts, execution time monitoring, and so on), there are no APIs to aid in the
7478 tuning of a realtime application’s behavior (**timespec** in timeouts, length of message queues,
7479 duration of driver interrupt service routine, and so on). The tuning of an application needs a
7480 means of recording timestamped important trace events during execution in order to analyze
7481 offline, and eventually, to tune some realtime features (redesign the system with less

7482 functionalities, readjust timeouts, redesign driver interrupts, and so on).

7483 Detailed Objectives

7484 Objectives were defined to build the trace interface and are kept for historical interest. Although
7485 some objectives are not fully respected in this trace interface, the concept of the POSIX trace
7486 interface assumes the following points:

- 7487 1. It shall be possible to trace both system and user trace events concurrently.
- 7488 2. It must be possible to trace per-process trace events and also to trace system trace events
7489 which are unrelated to any particular process. A per-process trace event is either user-
7490 initiated or system-initiated.
- 7491 3. It must be possible to control tracing on a per process basis from either inside or outside
7492 the process.
- 7493 4. It must be possible to control tracing on a per-thread basis from inside the enclosing
7494 process.
- 7495 5. Trace points shall be controllable by trace event type ID from inside and outside of the
7496 process. Multiple trace points can have the same trace event type ID, and will be controlled
7497 jointly.
- 7498 6. Recording of trace events is dependent on both trace event type ID and the
7499 process/thread. Both must be enabled in order to record trace events. System trace events
7500 may or may not be handled differently.
- 7501 7. The API shall not mandate the ability to control tracing for more than one process at the
7502 same time.
- 7503 8. There is no objective for trace control on anything bigger than a process; for example,
7504 group or session.
- 7505 9. Trace propagation and control:
 - 7506 a. Trace propagation across fork is optional; the default is to not trace a child process.
 - 7507 b. Trace control shall span *thread_create* operations; that is, if a process is being traced,
7508 any thread will be traced as well if this thread allows tracing. The default is to allow
7509 tracing.
- 7510 10. Trace control shall not span *exec* or *spawn* operations.
- 7511 11. A triggering API is not required. The triggering API is the ability to command or stop
7512 tracing based on the occurrence of specific trace event other than a `POSIX_TRACE_START`
7513 trace event or a `POSIX_TRACE_STOP` trace event.
- 7514 12. Trace log entries shall have timestamps of implementation-defined resolution.
7515 Implementations are exhorted to support at least microsecond resolution. When a trace log
7516 entry is retrieved, it shall have timestamp, PC address, PID, and TID of the entity that
7517 generated the trace event.
- 7518 13. Independently developed code should be able to use trace facilities without coordination
7519 and without conflict.
- 7520 14. Even if the trace points in the trace calls are not unique, the trace log entries (after any
7521 processing) shall be uniquely identified as to trace point.
- 7522 15. There shall be a standard API to read the trace stream.

- 7523 16. The format of the trace stream and the trace log is opaque and unspecified.
- 7524 17. It shall be possible to read a completed trace, if recorded on some suitable non-volatile
7525 storage, even subsequent to a power cycle or subsequent cold boot of the system.
- 7526 18. Support of analysis of a trace log while it is being formed is implementation-defined.
- 7527 19. The API shall allow the application to write trace stream identification information into the
7528 trace stream and to be able to retrieve it, without it being overwritten by trace entries, even
7529 if the trace stream is full.
- 7530 20. It must be possible to specify the destination of trace data produced by trace events.
- 7531 21. It must be possible to have different trace streams, and for the tracing enabled by one trace
7532 stream to be completely independent of the tracing of another trace stream.
- 7533 22. It must be possible to trace events from threads in different CPUs.
- 7534 23. The API shall support one or more trace streams per-system, and one or more trace
7535 streams per-process, up to an implementation-defined set of per-system and per-process
7536 maximums.
- 7537 24. It shall be possible to determine the order in which the trace events happened, without
7538 necessarily depending on the clock, up to an implementation-defined time resolution.
- 7539 25. For performance reasons, the trace event point call(s) shall be implementable as a macro
7540 (see the ISO POSIX-1: 1996 standard, 1.3.4, Statement 2).
- 7541 26. IEEE Std 1003.1-200x must not define the trace points which a conforming system must
7542 implement, except for trace points used in the control of tracing.
- 7543 27. The APIs shall be thread-safe, and trace points should be lock-free (that is shall not require
7544 a lock to gain exclusive access to some resource).
- 7545 28. The user-provided information associated with a trace event is variable-sized, up to some
7546 maximum size.
- 7547 29. Bounds on record and trace stream sizes:
- 7548 a. The API must permit the application to declare the upper bounds on the length of an
7549 application data record. The system shall return the limit it used. The limit used may
7550 be smaller than requested.
- 7551 b. The API must permit the application to declare the upper bounds on the size of trace
7552 streams. The system shall return the limit it used. The limit used may be different,
7553 either larger or smaller, than requested.
- 7554 30. The API must be able to pass any fundamental data type, and a structured data type
7555 composed only of fundamental types. The API must be able to pass data by reference,
7556 given only as an address and a length. Fundamental types are the POSIX.1 types (see the
7557 `<sys/types.h>` header) plus those defined in the ISO C standard.
- 7558 31. The API shall apply the POSIX notions of ownership and permission to recorded trace
7559 data, corresponding to the sources of that data.

7560 **Comments on Objectives**

7561 **Note:** In the following comments, numbers in square brackets refer to the above objectives.

7562 It is necessary to be able to obtain a trace stream for a complete activity. This means we need to
7563 be able to trace both application and system trace events. A per-process trace event is either
7564 user-initiated, like the *write()* POSIX call, or system-initiated, like a timer expiration. We also
7565 need to be able to trace an entire process's activity even when it has threads in multiple CPUs. To
7566 avoid excess trace activity, it is necessary to be able to control tracing on a trace event type basis.
7567 [Objectives 1,2,5,22]

7568 We need to be able to control tracing on a per-process basis, both from inside and outside the
7569 process; that is, a process can start a trace activity on itself or any other process. We also see the
7570 need to allow the definition of a maximum number trace streams per system.
7571 [Objectives 3,23]

7572 From within a process, it is necessary to be able to control tracing on a per-thread basis. This
7573 provides an additional filtering capability to keep the amount of traced data to a minimum. It
7574 also allows for less ambiguity as to the origin of trace events. It is recognized that thread-level
7575 control is only valid from within the process itself. It is also desirable to know the maximum
7576 number of trace streams per process that can be started. We do not want the API to require
7577 thread synchronization or to mandate priority inversions that would cause the thread to block.
7578 However, the API must be thread-safe.
7579 [Objectives 4,23,24,27]

7580 We see no objective to control tracing on anything larger than a process; for example, a group or
7581 session. Also, the ability to start or stop a trace activity on multiple processes atomically may be
7582 very difficult or cumbersome in some implementations.
7583 [Objectives 6,8]

7584 It is also necessary to be able to control tracing by trace event type identifier, sometimes called a
7585 trace hook ID. However, there is no mandated set of system trace events, since such trace points
7586 are very system-dependent. The API must not require from the operating system facilities that
7587 are not standard (POSIX).
7588 [Objectives 6,26]

7589 Trace control must span *fork()* and *pthread_create()*. If not, there will be no way to ensure that a
7590 program's activity is entirely traced. The newly forked child would not be able to turn on its
7591 tracing until after it obtained control after the fork, and trace control externally would be even
7592 more problematic.
7593 [Objective 9]

7594 Since *exec()* and *spawn()* represent a complete change in the execution of a task (a new
7595 program), trace control need not persist over an *exec()* or *spawn()*.
7596 [Objective 10]

7597 Where trace activities are started on multiple processes, these trace activities should not interfere
7598 with each other.
7599 [Objective 21]

7600 There is no need for a triggering objective, primarily for performance reasons; see also Section
7601 B.2.11.8 (on page 3491), rationale on triggering.
7602 [Objective 11]

7603 It must be possible to determine the origin of each traced event. We need the process and thread
7604 identifiers for each trace event. We also saw the need for a user-specifiable origin, but felt this
7605 would create too much overhead.
7606 [Objectives 12,14]

7607 We must allow for trace points to come embedded in software components from several
7608 different sources and vendors without requiring coordination.
7609 [Objective 13]

7610 We need to be able to uniquely identify trace points that may have the same trace stream
7611 identifier. We only need to be able to do this when a trace report is produced.
7612 [Objectives 12,14]

7613 Tracing is a very performance-sensitive activity, and will therefore likely be implemented at a
7614 low level within the system. Hence the interface shall not mandate any particular buffering or
7615 storage method. Therefore, we will need a standard API to read a trace stream. Also the interface
7616 shall not mandate the format of the trace data, and the interface shall not assume a trace storage
7617 method. Due to the possibility of a monolithic kernel and the possible presence of multiple
7618 processes capable of running trace activities, the two kinds of trace events may be stored in two
7619 separate streams for performance reasons. A mandatory dump mechanism, common in some
7620 existing practice, has been avoided to allow the implementation of this set of functions on small
7621 realtime profiles for which the concept of a file system is not defined. The trace API calls should
7622 be implemented as macros.
7623 [Objectives 15,16,25,30]

7624 Since a trace facility is a valuable service tool, the output (or log) of a completed trace stream
7625 that is written to permanent storage must be readable on other systems of the type that
7626 produced the trace log. Note that there is no objective to be able to interpret a trace log that was
7627 not successfully completed.
7628 [Objectives 17,18,19]

7629 For trace streams written to permanent storage, a way to specify the destination of the trace
7630 stream is needed.
7631 [Objective 20]

7632 We need to be able to depend on the ordering of trace events up to some implementation- |
7633 defined time interval. For example, we need to know the time period which, if trace events are |
7634 closer together, their ordering is unspecified. Events that occur within an interval smaller than |
7635 this resolution may or may not be read back in the correct order. |
7636 [Objective 24]

7637 The application should be able to know how much data can be traced. When trace event types
7638 can be filtered, the application should be able to specify the approximate maximum amount of
7639 data that will be traced in a trace event so resources can be more efficiently allocated.
7640 [Objectives 28,29]

7641 Users should not be able to trace data to which they would not normally have access to. System
7642 trace events corresponding to a process/thread should be associated with the ownership of that
7643 process/thread.
7644 [Objective 31] |

7645 B.2.11.2 Trace Model

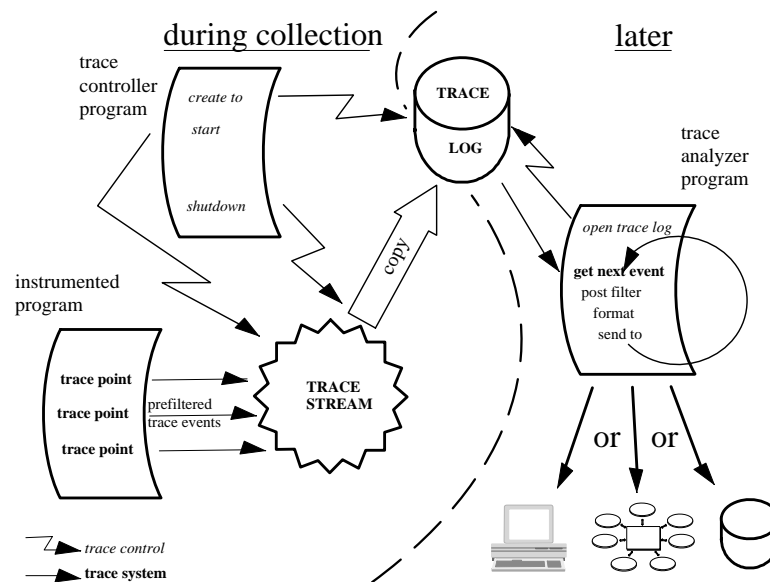
7646 **Introduction**

7647 The model is based on two base entities: the “Trace Stream” and the “Trace Log” , and a
 7648 recorded unit called the “Trace Event”. The possibility of using Trace Streams and Trace Logs
 7649 separately gives us two use dimensions and solves both the performance issue and the full-
 7650 information system issue. In the case of a trace stream without log, specific information,
 7651 although reduced in quantity, is required to be registered, in a possibly small realtime system,
 7652 with as little overhead as possible. The Trace Log option has been added for small realtime
 7653 systems. In the case of a trace stream with log, considerable complex application-specific
 7654 information needs to be collected.

7655 **Trace Model Description**

7656 The trace model can be examined for three different subfunctions: Application Instrumentation,
 7657 Trace Operation Control, and Trace Analysis.

7658



7659 **Figure B-2** Trace System Overview: for Offline Analysis

7660 Each of these subfunctions requires specific characteristics of the trace mechanism API.

7661 • Application Instrumentation

7662 When instrumenting an application, the programmer has no concern about the future
 7663 utilization of the trace events in trace stream or trace log, the full policy of trace stream, or
 7664 the eventual pre-filtering of trace events. But he is concerned about the correct determination
 7665 of specific trace event type identifier, regardless of how many independent libraries are used
 7666 in the same user application; see Figure B-2 and Figure B-3 (on page 3475).

7667 This trace API shall provide the necessary operations to accomplish this subfunction. This is
7668 done by providing functions to associate a programmer-defined name with an
7669 implementation-defined trace event type identifier; see the *posix_trace_eventid_open()*
7670 function), and to send this trace event into a potential trace stream (see the
7671 *posix_trace_event()* function).

7672 • Trace Operation Control

7673 When controlling the recording of trace events in a trace stream, the programmer is
7674 concerned with the correct initialization of the trace mechanism (that is, the sizing of the
7675 trace stream), the correct retention of trace events in a permanent storage, the correct
7676 dynamic recording of trace events, and so on.

7677 This trace API shall provide the necessary material to permit this efficiently. This is done by
7678 providing functions to initialize a new trace stream, and optionally a trace log:

- 7679 — Trace Stream Attributes Object Initialization (see *posix_trace_attr_init()*)
- 7680 — Functions to Retrieve or Set Information About a Trace Stream (see
7681 *posix_trace_attr_getgenversion()*)
- 7682 — Functions to Retrieve or Set the Behavior of a Trace Stream (see
7683 *posix_trace_attr_getinherited()*)
- 7684 — Functions to Retrieve or Set Trace Stream Size Attributes (see
7685 *posix_trace_attr_getmaxusevents()*)
- 7686 — Trace Stream Initialization, Flush, and Shutdown from a Process (see *posix_trace_create()*)
- 7687 — Clear Trace Stream and Trace Log (see *posix_trace_clear()*)

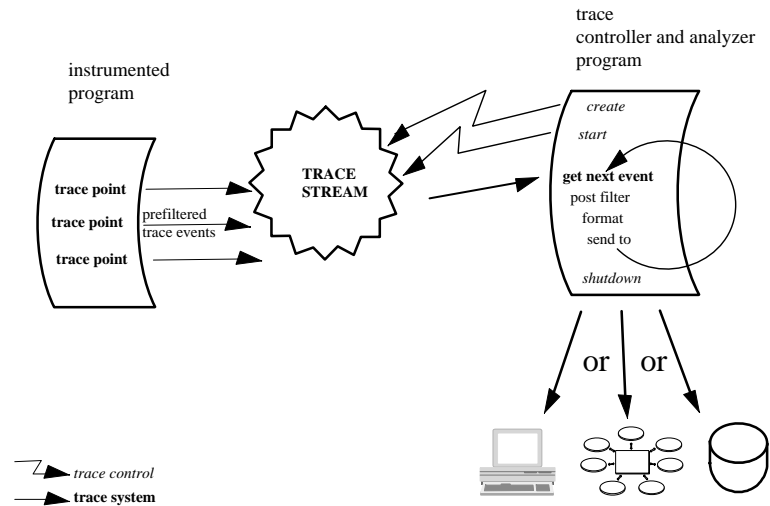
7688 To select the trace event types that are to be traced:

- 7689 — Manipulate Trace Event Type Identifier (see *posix_trace_trid_eventid_open()*)
- 7690 — Iterate over a Mapping of Trace Event Type (see *posix_trace_eventtypelist_getnext_id()*)
- 7691 — Manipulate Trace Event Type Sets (see *posix_trace_eventset_empty()*)
- 7692 — Set Filter of an Initialized Trace Stream (see *posix_trace_set_filter()*)

7693 To control the execution of an active trace stream:

- 7694 — Trace Start and Stop (see *posix_trace_start()*)
- 7695 — Functions to Retrieve the Trace Attributes or Trace Statuses (see *posix_trace_get_attr()*)

7696



7697

Figure B-3 Trace System Overview: for Online Analysis

7698

- Trace Analysis

7699

7700

7701

Once correctly recorded, on permanent storage or not, an ultimate activity consists of the analysis of the recorded information. If the recorded data is on permanent storage, a specific open operation is required to associate a trace stream to a trace log.

7702

7703

7704

7705

7706

7707

The first intent of the group was to request the presence of a system identification structure in the trace stream attribute. This was, for the application, to allow some portable way to process the recorded information. However, there is no requirement that the **utsname** structure, on which this system identification was based, be portable from one machine to another, so the contents of the attribute cannot be interpreted correctly by an application conforming to IEEE Std 1003.1-200x.

7708

7709

7710

7711

7712

Draft 6 incorporates this modification and requests that some unspecified information be recorded in the trace log in order to fail opening it if the analysis process and the controller process were running in different types of machine, but does not request that this information be accessible to the application. This modification has implied a modification in the *posix_trace_open()* function error code returns.

7713

This trace API shall provide functions to:

7714

- Extract trace stream identification attributes (see *posix_trace_attr_getgenversion()*)

7715

- Extract trace stream behavior attributes (see *posix_trace_attr_getinherited()*)

7716

7717

- Extract trace event, stream, and log size attributes (see *posix_trace_attr_getmaxusereventsize()*)

7718

- Look up trace event type names (see *posix_trace_eventid_get_name()*)

- 7719 — Iterate over trace event type identifiers (see *posix_trace_eventtypelist_getnext_id()*)
- 7720 — Open, rewind, and close a trace log (see *posix_trace_open()*)
- 7721 — Read trace stream attributes and status (see *posix_trace_get_attr()*)
- 7722 — Read trace events (see *posix_trace_getnext_event()*)

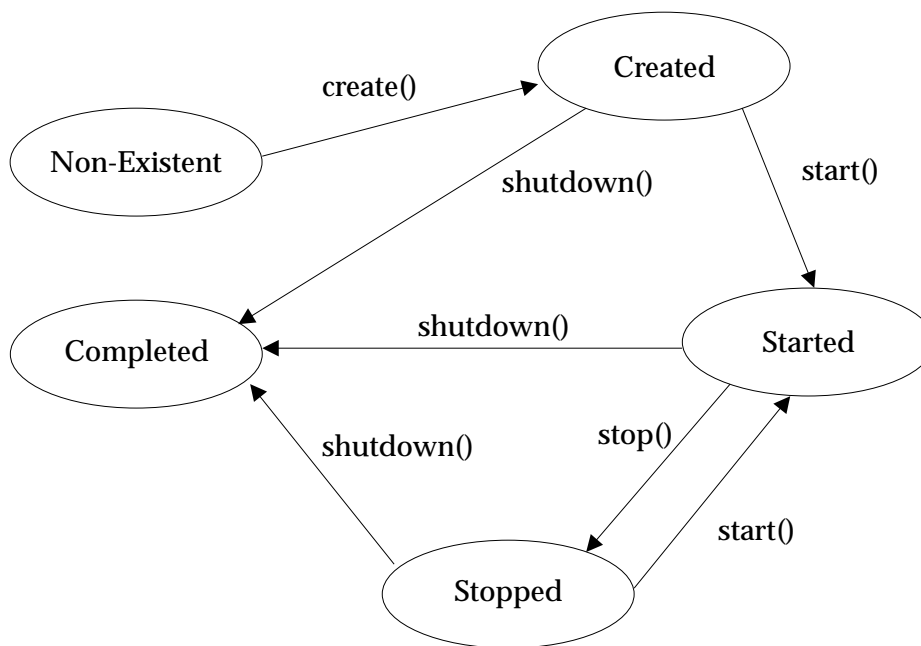
7723 Due to the following two reasons:

- 7724 1. The requirement that the trace system must not add unacceptable overhead to the traced
7725 process and so that the trace event point execution must be fast
- 7726 2. The traced application does not care about tracing errors

7727 the trace system cannot return any internal error to the application. Internal error conditions can
7728 range from unrecoverable errors that will force the active trace stream to abort, to small errors
7729 that can affect the quality of tracing without aborting the trace stream. The group decided to
7730 define a system trace event to report to the analysis process such internal errors. It is not the
7731 intention of IEEE Std 1003.1-200x to require an implementation to report an internal error that
7732 corrupts or terminates tracing operation. The implementor is free to decide which internal
7733 documented errors, if any, the trace system is able to report.

7734 States of a Trace Stream

7735



7736

Figure B-4 Trace System Overview: States of a Trace Stream

7737 Figure B-4 shows the different states an active trace stream passes through. After the
7738 *posix_trace_create()* function call, a trace stream becomes CREATED and a trace stream is
7739 associated for the future collection of trace events. The status of the trace stream is
7740 POSIX_TRACE_SUSPENDED. The state becomes STARTED after a call to the *posix_trace_start()*
7741 function, and the status becomes POSIX_TRACE_RUNNING. In this state, all trace events that
7742 are not filtered out shall be stored into the trace stream. After a call to *posix_trace_stop()*,
7743 the trace stream becomes STOPPED (and the status POSIX_TRACE_SUSPENDED). In this state, no

7744 new trace events will be recorded in the trace stream, but previously recorded trace events may
7745 continue to be read.

7746 After a call to *posix_trace_shutdown()*, the trace stream is in the state COMPLETED. The trace
7747 stream no longer exists but, if the Trace Log option is supported, all the information contained in
7748 it has been logged. If a log object has not been associated with the trace stream at the creation, it
7749 is the responsibility of the trace controller process to not shut the trace stream down while trace
7750 events remain to be read in the stream.

7751 **Tracing All Processes**

7752 Some implementations have a tracing subsystem with the ability to trace all processes. This is
7753 useful to debug some types of device drivers such as those for ATM or X25 adapters. These types
7754 of adapters are used by several independent processes, that are not issued from the same
7755 process.

7756 The POSIX trace interface does not define any constant or option to create a trace stream tracing
7757 all processes. But the POSIX trace interface does not prevent this type of implementation and the
7758 implementor is free to add this capability. Nevertheless, the POSIX trace interface allows to trace
7759 all the system trace events and all the processes issued from the same process.

7760 If such a tracing system capability has to be implemented, when a trace stream is created, it is
7761 recommended that a constant named POSIX_TRACE_ALLPROC be used instead of the process
7762 identifier in the argument of the function *posix_trace_create()* or *posix_trace_create_withlog()*. A
7763 possible value for POSIX_TRACE_ALLPROC may be -1 instead of a real process identifier.

7764 The implementor has to be aware that there is some impact on the tracing behavior as defined in
7765 the POSIX trace interface. For example:

- 7766 • If the default value for the inheritance attribute is to set to
7767 POSIX_TRACE_CLOSE_FOR_CHILD, the implementation has to stop tracing for the child
7768 process.
- 7769 • The trace controller which is creating this type of trace stream must have the appropriate
7770 privilege to trace all the processes.

7771 **Trace Storage**

7772 The model is based on two types of trace events: system trace events and user-defined trace
7773 events. The internal representation of trace events is implementation-defined, and so the
7774 implementor is free to choose the more suitable, practical, and efficient way to design the
7775 internal management of trace events. For the timestamping operation, the model does not
7776 impose the CLOCK_REALTIME or any other clock. The buffering allocation and operation
7777 follow the same principle. The implementor is free to use one or more buffers to record trace
7778 events; the interface assumes only a logical trace stream of sequentially recorded trace events.
7779 Regarding flushing of trace events, the interface allows the definition of a trace log object which
7780 typically can be a file. But the group was also aware of defining functions to permit the use of
7781 this interface in small realtime systems, which may not have general file system capabilities. For
7782 instance, the three functions *posix_trace_getnext_event()* (blocking),
7783 *posix_trace_timedgetnext_event()* (blocking with timeout), and *posix_trace_trygetnext_event()*
7784 (non-blocking) are proposed to read the recorded trace events.

7785 The policy to be used when the trace stream becomes full also relies on common practice:

- 7786 • For an active trace stream, the POSIX_TRACE_LOOP trace stream policy permits automatic
7787 overrun (overwrite of oldest trace events) while waiting for some user-defined condition to
7788 cause tracing to stop. By contrast, the POSIX_TRACE_UNTIL_FULL trace stream policy

7789 requires the system to stop tracing when the trace stream is full. However, if the trace stream
 7790 that is full is at least partially emptied by a call to the *posix_trace_flush()* function or by calls
 7791 to *posix_trace_getnext_event()* function, the trace system will automatically resume tracing.

7792 If the Trace Log option is supported the operation of the POSIX_TRACE_FLUSH policy is an
 7793 extension of the POSIX_TRACE_UNTIL_FULL policy. The automatic free operation (by
 7794 flushing to the associated trace log) is added.

7795 • If a log is associated with the trace stream and this log is a regular file, these policies also
 7796 apply for the log. One more policy, POSIX_TRACE_APPEND, is defined to allow indefinite
 7797 extension of the log. Since the log destination can be any device or pseudo-device, the
 7798 implementation may not be able to manipulate the destination as required by
 7799 IEEE Std 1003.1-200x. For this reason, the behavior of the log full policy may be unspecified
 7800 depending of the trace log type.

7801 The current trace interface does not define a service to preallocate space for a trace log file,
 7802 because this space can be preallocated by means of a call to the *posix_fallocate()* function. This
 7803 function could be called after the file has been opened, but before the trace stream is created.
 7804 The *posix_fallocate()* function ensures that any required storage for regular file data is
 7805 allocated on the file system storage media. If *posix_fallocate()* returns successfully,
 7806 subsequent writes to the specified file data shall not fail due to the lack of free space on the
 7807 file system storage media. Besides trace events, a trace stream also includes trace attributes
 7808 and the mapping from trace event names to trace event type identifiers. The implementor is
 7809 free to choose how to store the trace attributes and the trace event type map, but must ensure
 7810 that this information is not lost when a trace stream overrun occurs.

7811 B.2.11.3 Trace Programming Examples

7812 Several programming examples are presented to show the code of the different possible
 7813 subfunctions using a trace subsystem. All these programs need to include the `<trace.h>` header.
 7814 In the examples shown, error checking is omitted for more simplicity.

7815 Trace Operation Control

7816 These examples show the creation of a trace stream for another process; one which is already
 7817 trace instrumented. All the default trace stream attributes are used to simplify programming in
 7818 the first example. The second example shows more possibilities.

7819 First Example

```
7820 /* Caution. Error checks omitted */
7821 {
7822     trace_attr_t attr;
7823     pid_t pid = traced_process_pid;
7824     int fd;
7825     trace_id_t trid;
7826     - - - - -
7827     /* Initialize trace stream attributes */
7828     posix_trace_attr_init(&attr);
7829     /* Open a trace log */
7830     fd=open("/tmp/mytracelog",...);
7831     /*
7832      * Create a new trace associated with a log
7833      * and with default attributes
7834      */
```

```

7835     posix_trace_create_withlog(pid, &attr, fd, &trid);
7836     /* Trace attribute structure can now be destroyed */
7837     posix_trace_attr_destroy(&attr);
7838     /* Start of trace event recording */
7839     posix_trace_start(trid);
7840     - - - - -
7841     - - - - -
7842     /* Duration of tracing */
7843     - - - - -
7844     - - - - -
7845     /* Stop and shutdown of trace activity */
7846     posix_trace_shutdown(trid);
7847     - - - - -
7848 }

```

7849 **Second Example**

7850 Between the initialization of the trace stream attributes and the creation of the trace stream,
7851 these trace stream attributes may be modified; see **Trace Stream Attribute Manipulation** (on
7852 page 3483) for specific programming example. Between the creation and the start of the trace
7853 stream, the event filter may be set; after the trace stream is started, the event filter may be
7854 changed. The setting of an event set and the change of a filter is shown in **Create a Trace Event**
7855 **Type Set and Change the Trace Event Type Filter** (on page 3483).

```

7856 /* Caution. Error checks omitted */
7857 {
7858     trace_attr_t attr;
7859     pid_t pid = traced_process_pid;
7860     int fd;
7861     trace_id_t trid;
7862     - - - - -
7863     /* Initialize trace stream attributes */
7864     posix_trace_attr_init(&attr);
7865     /* Attr default may be changed at this place; see example */
7866     - - - - -
7867     /* Create and open a trace log with R/W user access */
7868     fd=open("/tmp/mytracelog",O_WRONLY|O_CREAT,S_IRUSR|S_IWUSR);
7869     /* Create a new trace associated with a log */
7870     posix_trace_create_withlog(pid, &attr, fd, &trid);
7871     /*
7872      * If the Trace Filter option is supported
7873      * trace event type filter default may be changed at this place;
7874      * see example about changing the trace event type filter
7875      */
7876     posix_trace_start(trid);
7877     - - - - -
7878     /*
7879      * If you have an uninteresting part of the application
7880      * you can stop temporarily.
7881      *
7882      * posix_trace_stop(trid);
7883      * - - - - -

```

```

7884         * - - - - -
7885         * posix_trace_start(trid);
7886         */
7887     - - - - -
7888     /*
7889     * If the Trace Filter option is supported
7890     * the current trace event type filter can be changed
7891     * at any time (see example about how to set
7892     * a trace event type filter
7893     */
7894     - - - - -
7895     /* Stop the recording of trace events */
7896     posix_trace_stop(trid);
7897     /* Shutdown the trace stream */
7898     posix_trace_shutdown(trid);
7899     /*
7900     * Destroy trace stream attributes; attr structure may have
7901     * been used during tracing to fetch the attributes
7902     */
7903     posix_trace_attr_destroy(&attr);
7904     - - - - -
7905 }

```

7906 **Application Instrumentation**

7907 This example shows an instrumented application. The code is included in a block of instructions,
7908 perhaps a function from a library. Possibly in an initialization part of the instrumented
7909 application, two user trace event names are mapped to two trace event type identifiers
7910 (function *posix_trace_eventid_open()*). Then two trace points are programmed.

```

7911 /* Caution. Error checks omitted */
7912 {
7913     trace_event_id_t eventid1, eventid2;
7914     - - - - -
7915     /* Initialization of two trace event type ids */
7916     posix_trace_eventid_open("my_first_event",&eventid1);
7917     posix_trace_eventid_open("my_second_event",&eventid2);
7918     - - - - -
7919     - - - - -
7920     - - - - -
7921     /* Trace point */
7922     posix_trace_event(eventid1,NULL,0);
7923     - - - - -
7924     /* Trace point */
7925     posix_trace_event(eventid2,NULL,0);
7926     - - - - -
7927 }

```

7928 **Trace Analyzer**

7929 This example shows the manipulation of a trace log resulting from the dumping of a completed
 7930 trace stream. All the default attributes are used to simplify programming, and data associated
 7931 with a trace event are not shown in the first example. The second example shows more
 7932 possibilities.

7933 **First Example**

```

7934        /* Caution. Error checks omitted */
7935        {
7936            int fd;
7937            trace_id_t trid;
7938            posix_trace_event_info trace_event;
7939            char trace_event_name[TRACE_EVENT_NAME_MAX];
7940            int return_value;
7941            size_t returndatasize;
7942            int lost_event_number;
7943            - - - - -
7944            /* Open an existing trace log */
7945            fd=open("/tmp/tracelog", O_RDONLY);
7946            /* Open a trace stream on the open log */
7947            posix_trace_open(fd, &trid);
7948            /* Read a trace event */
7949            posix_trace_getnext_event(trid, &trace_event,
7950                                    NULL, 0, &returndatasize,&return_value);
7951            /* Read and print all trace event names out in a loop */
7952            while (return_value == NULL)
7953            {
7954                /*
7955                * Get the name of the trace event associated
7956                * with trid trace ID
7957                */
7958                posix_trace_eventid_get_name(trid, trace_event.event_id,
7959                                            trace_event_name);
7960                /* Print the trace event name out */
7961                printf("%s\n",trace_event_name);
7962                /* Read a trace event */
7963                posix_trace_getnext_event(trid, &trace_event,
7964                                            NULL, 0, &returndatasize,&return_value);
7965            }
7966            /* Close the trace stream */
7967            posix_trace_close(trid);
7968            /* Close the trace log */
7969            close(fd);
7970        }

```


7971 **Second Example**

7972 The complete example includes the two other examples in **Retrieve Information from a Trace Log** (on page 3484) and in **Retrieve the List of Trace Event Types Used in a Trace Log** (on page 3485). For example, the *maxdatasize* variable is set in **Retrieve the List of Trace Event Types Used in a Trace Log** (on page 3485).

```

7976 /* Caution. Error checks omitted */
7977 {
7978     int fd;
7979     trace_id_t trid;
7980     posix_trace_event_info trace_event;
7981     char trace_event_name[TRACE_EVENT_NAME_MAX];
7982     char * data;
7983     size_t maxdatasize=1024, returndatasize;
7984     int return_value;
7985     - - - - -
7986     /* Open an existing trace log */
7987     fd=open("/tmp/tracelog", O_RDONLY);
7988     /* Open a trace stream on the open log */
7989     posix_trace_open( fd, &trid);
7990     /*
7991      * Retrieve information about the trace stream which
7992      * was dumped in this trace log (see example)
7993      */
7994     - - - - -
7995     /* Allocate a buffer for trace event data */
7996     data=(char *)malloc(maxdatasize);
7997     /*
7998      * Retrieve the list of trace event used in this
7999      * trace log (see example)
8000      */
8001     - - - - -
8002     /* Read and print all trace event names and data out in a loop */
8003     while (1)
8004     {
8005     posix_trace_getnext_event(trid, &trace_event,
8006         data, maxdatasize, &returndatasize,&return_value);
8007         if (return_value != NULL) break;
8008         /*
8009          * Get the name of the trace event type associated
8010          * with trid trace ID
8011          */
8012         posix_trace_eventid_get_name(trid, trace_event.event_id,
8013             trace_event_name);
8014         {
8015         int i;
8016
8017         /* Print the trace event name out */
8018         printf("%s: ", trace_event_name);
8019         /* Print the trace event data out */
8020         for (i=0; i<returndatasize, i++) printf("%02.2X",

```

```

8020         (unsigned char)data[i]);
8021     printf("\n");
8022     }
8023 }
8024     /* Close the trace stream */
8025     posix_trace_close(trid);
8026     /* The buffer data is deallocated */
8027     free(data);
8028     /* Now the file can be closed */
8029     close(fd);
8030 }

```

8031 **Several Programming Manipulations**

8032 The following examples show some typical sets of operations needed in some contexts.

8033 **Trace Stream Attribute Manipulation**

8034 This example shows the manipulation of a trace stream attribute object in order to change the
8035 default value provided by a previous *posix_trace_attr_init()* call.

```

8036     /* Caution. Error checks omitted */
8037     {
8038         trace_attr_t attr;
8039         size_t logsize=100000;
8040         - - - - -
8041         /* Initialize trace stream attributes */
8042         posix_trace_attr_init(&attr);
8043         /* Set the trace name in the attributes structure */
8044         posix_trace_attr_setname(&attr, "my_trace");
8045         /* Set the trace full policy */
8046         posix_trace_attr_setstreamfullpolicy(&attr, POSIX_TRACE_LOOP);
8047         /* Set the trace log size */
8048         posix_trace_attr_setlogsize(&attr, logsize);
8049         - - - - -
8050     }

```

8051 **Create a Trace Event Type Set and Change the Trace Event Type Filter**

8052 This example is valid only if the Trace Event Filter option is supported. This example shows the
8053 manipulation of a trace event type set in order to change the trace event type filter for an existing
8054 active trace stream, which may be just-created, running, or suspended. Some sets of trace event
8055 types are well-known, such as the set of trace event types not associated with a process, some
8056 trace event types are just-built trace event types for this trace stream; one trace event type is the
8057 predefined trace event error type which is deleted from the trace event type set.

```

8058     /* Caution. Error checks omitted */
8059     {
8060         trace_id_t trid = existing_trace;
8061         trace_event_set_t set;
8062         trace_event_id_t trace_event1, trace_event2;
8063         - - - - -
8064         /* Initialize to an empty set of trace event types */

```

```

8065      /* (not strictly required because posix_trace_event_set_fill() */ |
8066      /* will ignore the prior contents of the event set.) */ |
8067      posix_trace_eventset_emptyset(&set); |
8068      /* |
8069      * Fill the set with all system trace events |
8070      * not associated with a process |
8071      */ |
8072      posix_trace_eventset_fill(&set, POSIX_TRACE_WOPID_EVENTS); |
8073      /* |
8074      * Get the trace event type identifier of the known trace event name |
8075      * my_first_event for the trid trace stream |
8076      */ |
8077      posix_trace_trid_eventid_open(trid, "my_first_event", &trace_event1); |
8078      /* Add the set with this trace event type identifier */ |
8079      posix_trace_eventset_add_event(trace_event1, &set); |
8080      /* |
8081      * Get the trace event type identifier of the known trace event name |
8082      * my_second_event for the trid trace stream |
8083      */ |
8084      posix_trace_trid_eventid_open(trid, "my_second_event", &trace_event2); |
8085      /* Add the set with this trace event type identifier */ |
8086      posix_trace_eventset_add_event(trace_event2, &set); |
8087      - - - - - |
8088      /* Delete the system trace event POSIX_TRACE_ERROR from the set */ |
8089      posix_trace_eventset_del_event(POSIX_TRACE_ERROR, &set); |
8090      - - - - - |
8091      /* Modify the trace stream filter making it equal to the new set */ |
8092      posix_trace_set_filter(trid, &set, POSIX_TRACE_SET_EVENTSET); |
8093      - - - - - |
8094      /* |
8095      * Now trace_event1, trace_event2, and all system trace event types |
8096      * not associated with a process, except for the POSIX_TRACE_ERROR |
8097      * system trace event type, are filtered out of (not recorded in) the |
8098      * existing trace stream. |
8099      */ |
8100      }

```

8101 **Retrieve Information from a Trace Log**

8102 This example shows how to extract information from a trace log, the dump of a trace stream.
8103 This code:

```

8104      • Asks if the trace stream has lost trace events
8105      • Extracts the information about the version of the trace subsystem which generated this trace
8106      log
8107      • Retrieves the maximum size of trace event data; this may be used to dynamically allocate an
8108      array for extracting trace event data from the trace log without overflow
8109      /* Caution. Error checks omitted */
8110      {
8111          struct posix_trace_status_info statusinfo;

```

```

8112     trace_attr_t attr;
8113     trace_id_t trid = existing_trace;
8114     size_t maxdatasize;
8115     char genversion[TRACE_NAME_MAX];
8116     - - - - -
8117     /* Get the trace stream status */
8118     posix_trace_get_status(trid, &statusinfo);
8119     /* Detect an overrun condition */
8120     if (statusinfo.posix_stream_overrun_status == POSIX_TRACE_OVERRUN)
8121         printf("trace events have been lost\n");

8122     /* Get attributes from the trid trace stream */
8123     posix_trace_get_attr(trid, &attr);
8124     /* Get the trace generation version from the attributes */
8125     posix_trace_attr_getgenversion(&attr, genversion);
8126     /* Print the trace generation version out */
8127     printf("Information about Trace Generator:%s\n",genversion);

8128     /* Get the trace event max data size from the attributes */
8129     posix_trace_attr_getmaxdatasize(&attr, &maxdatasize);
8130     /* Print the trace event max data size out */
8131     printf("Maximum size of associated data:%d\n",maxdatasize);
8132     /* Destroy the trace stream attributes */
8133     posix_trace_attr_destroy(&attr);
8134 }

```

8135 **Retrieve the List of Trace Event Types Used in a Trace Log**

8136 This example shows the retrieval of a trace stream's trace event type list. This operation may be
8137 very useful if you are interested only in tracking the type of trace events in a trace log.

```

8138 /* Caution. Error checks omitted */
8139 {
8140     trace_id_t trid = existing_trace;
8141     trace_event_id_t event_id;
8142     char event_name[TRACE_EVENT_NAME_MAX];
8143     int return_value;
8144     - - - - -

8145     /*
8146      * In a loop print all existing trace event names out
8147      * for the trid trace stream
8148      */
8149     while (1)
8150     {
8151         posix_trace_eventtypelist_getnext_id(trid, &event_id
8152             &return_value);
8153         if (return_value != NULL) break;
8154         /*
8155          * Get the name of the trace event associated
8156          * with trid trace ID
8157          */
8158         posix_trace_eventid_get_name(trid, event_id, event_name);
8159         /* Print the name out */

```

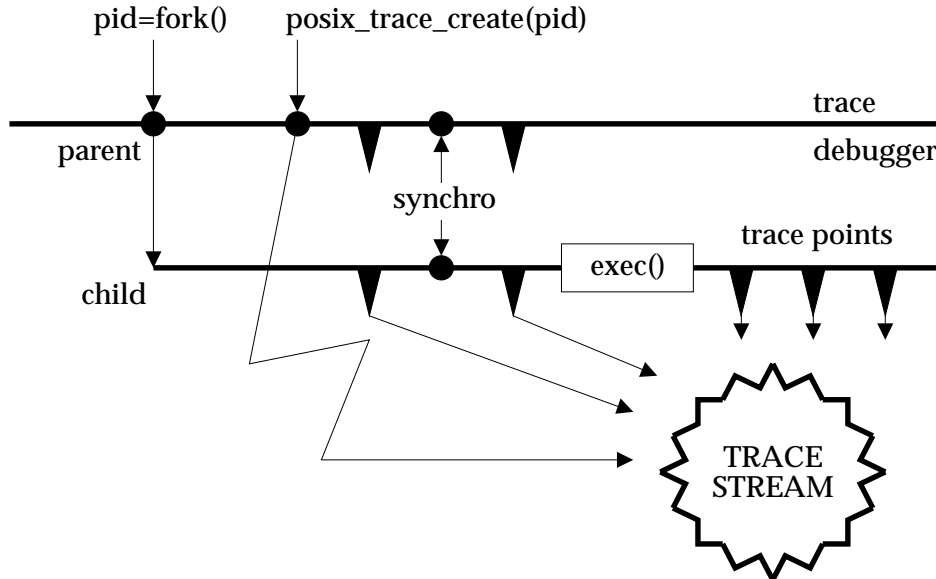
```

8160         printf("%s\n", event_name);
8161     }
8162 }

```

8163 B.2.11.4 Rationale on Trace for Debugging

8164



8165

Figure B-5 Trace Another Process

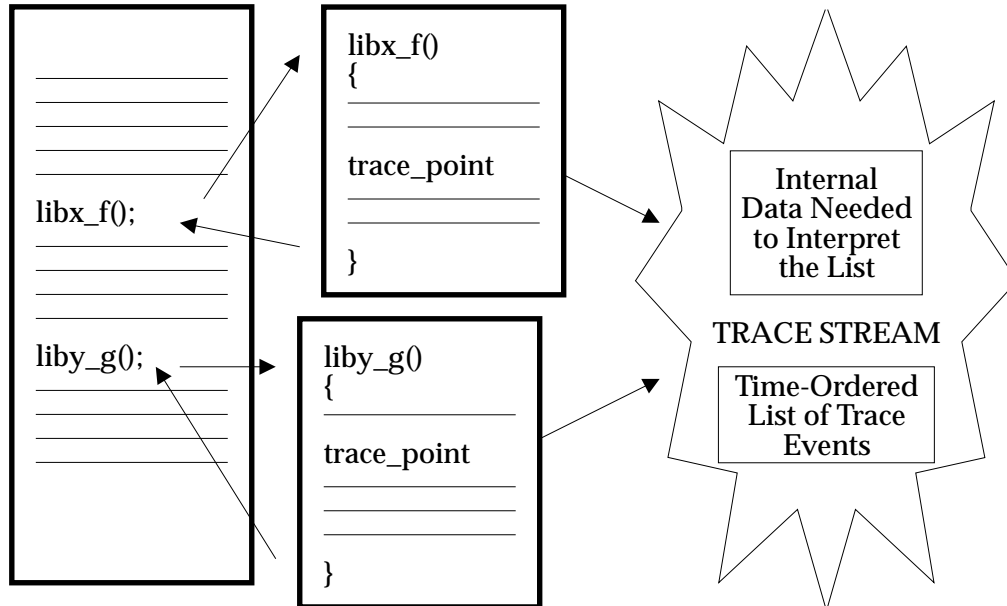
8166 Among the different possibilities offered by the trace interface defined in IEEE Std 1003.1-200x,
 8167 the debugging of an application is the most interesting one. Typical operations in the controlling
 8168 debugger process are to filter trace event types, to get trace events from the trace stream, to stop
 8169 the trace stream when the debugged process is executing uninteresting code, to start the trace
 8170 stream when some interesting point is reached, and so on. The interface defined in
 8171 IEEE Std 1003.1-200x should define all the necessary base functions to allow this dynamic debug
 8172 handling.

8173 Figure B-5 shows an example in which the trace stream is created after the call to the *fork()*
 8174 function. If the user does not want to lose trace events some synchronization mechanism
 8175 (represented in the figure) may be needed before calling the *exec()* function, to give the parent a
 8176 chance to create the trace stream before the child begins the execution of its trace points.

8177 B.2.11.5 Rationale on Trace Event Type Name Space

8178 At first, the working group was in favor of the representation of a trace event type by an integer
 8179 (*event_name*). It seems that existing practice shows the weakness of such a representation. The
 8180 collision of trace event types is the main problem that cannot be simply resolved using this sort
 8181 of representation. Suppose, for example, that a third party designs an instrumented library. The
 8182 user does not have the source of this library and wants to trace his application which uses in
 8183 some part the third-party library. There is no means for him to know what are the trace event
 8184 types used in the instrumented library so he has some chance of duplicating some of them and
 8185 thus to obtain a contaminated tracing of his application.

8186



8187

Figure B-6 Trace Name Space Overview: With Third-Party Library

8188 We have requirements to allow program images containing pieces from various vendors to be
 8189 traced without also requiring those or any other vendors to coordinate their uses of the trace
 8190 facility, and especially the naming of their various trace event types and trace point IDs. The
 8191 chosen solution is to provide a very large name space, large enough so that the individual
 8192 vendors can give their trace types and tracepoint IDs sufficiently long and descriptive names
 8193 making the occurrence of collisions quite unlikely. The probability of collision is thus made
 8194 sufficiently low so that the problem may, as a practical matter, be ignored. By requirement, the
 8195 consequence of collisions will be a slight ambiguity in the trace streams; tracing will continue in
 8196 spite of collisions and ambiguities. “The show must go on”. The *posix_prog_address* member of
 8197 the **posix_trace_event_info** structure is used to allow trace streams to be unambiguously
 8198 interpreted, despite the fact that trace event types and trace event names need not be unique.

8199 The *posix_trace_eventid_open()* function is required to allow the instrumented third-party library
 8200 to get a valid trace event type identifier for its trace event names. This operation is, somehow,
 8201 an allocation, and the group was aware of proposing some deallocation mechanism which the
 8202 instrumented application could use to recover the resources used by a trace event type identifier.
 8203 This would have given the instrumented application the benefit of being capable of reusing a
 8204 possible minimum set of trace event type identifiers, but also the inconvenience to have,
 8205 possibly in the same trace stream, one trace event type identifier identifying two different trace
 8206 event types. After some discussions the group decided to not define such a function which
 8207 would make this API thicker for little benefit, the user having always the possibility of adding
 8208 identification information in the data member of the trace event structure.

8209 The set of the trace event type identifiers the controlling process wants to filter out is initialized
 8210 in the trace mechanism using the function *posix_trace_set_filter()*, setting the arguments
 8211 according to the definitions explained in *posix_trace_set_filter()*. This operation can be done
 8212 statically (when the trace is in the STOPPED state) or dynamically (when the trace is in the
 8213 STARTED state). The preparation of the filter is normally done using the function defined in
 8214 *posix_trace_eventtypelist_getnext_id()* and eventually the function
 8215 *posix_trace_eventtypelist_rewind()* in order to know (before the recording) the list of the potential

8216 set of trace event types that can be recorded. In the case of an active trace stream, this list may
8217 not be exhaustive. Actually, the target process may not have yet called the function
8218 *posix_trace_eventid_open()*. But it is a common practice, for a controlling process, to prepare the
8219 filtering of a future trace stream before its start. Therefore the user must have a way to get the
8220 trace event type identifier corresponding to a well-known trace event name before its future
8221 association by the pre-cited function. This is done by calling the *posix_trace_trid_eventid_open()*
8222 function, given the trace stream identifier and the trace name, and described hereafter. Because
8223 this trace event type identifier is associated with a trace stream identifier, where a unique
8224 process has initialized two or more traces, the implementation is expected to return the same
8225 trace event type identifier for successive calls to *posix_trace_trid_eventid_open()* with different
8226 trace stream identifiers. The *posix_trace_eventid_get_name()* function is used by the controller
8227 process to identify, by the name, the trace event type returned by a call to the
8228 *posix_trace_eventtypelist_getnext_id()* function.

8229 Afterwards, the set of trace event types is constructed using the functions defined in
8230 *posix_trace_eventset_empty()*, *posix_trace_eventset_fill()*, *posix_trace_eventset_add()*, and
8231 *posix_trace_eventset_del()*.

8232 A set of functions is provided devoted to the manipulation of the trace event type identifier and
8233 names for an active trace stream. All these functions require the trace stream identifier argument
8234 as the first parameter. The opacity of the trace event type identifier implies that the user cannot
8235 associate directly its well-known trace event name with the system associated trace event type
8236 identifier.

8237 The *posix_trace_trid_eventid_open()* function allows the application to get the system trace event
8238 type identifier back from the system, given its well-known trace event name. This function is
8239 useful only when a controlling process needs to specify specific events to be filtered.

8240 The *posix_trace_eventid_get_name()* function allows the application to obtain a trace event name
8241 given its trace event type identifier. One possible use of this function is to identify the type of a
8242 trace event retrieved from the trace stream, and print it. The easiest way to implement this
8243 requirement, is to use a single trace event type map for all the processes whose maps are
8244 required to be identical. A more difficult way is to attempt to keep multiple maps identical at
8245 every call to *posix_trace_eventid_open()* and *posix_trace_trid_eventid_open()*.

8246 *B.2.11.6 Rationale on Trace Events Type Filtering*

8247 The most basic rationale for runtime and pre-registration filtering (selection/rejection) of trace
8248 event types is to prevent choking of the trace collection facility, and/or overloading of the
8249 computer system. Any worthwhile trace facility can bring even the largest computer to its
8250 knees. Otherwise, we would record everything, and filter after the fact; it would be much
8251 simpler, but impractical.

8252 To achieve debugging, measurement, or whatever the purpose of tracing, the filtering of trace
8253 event types is an important part of trace analysis. Due to the fact that the trace events are put
8254 into a trace stream and probably logged afterwards into a file, different levels of filtering—that
8255 is, rejection of trace event types—are possible.

8256 **Filtering of Trace Event Types Before Tracing**

8257 This function, represented by the *posix_trace_set_filter()* function in IEEE Std 1003.1-200x (see
 8258 *posix_trace_set_filter()*), selects, before or during tracing, the set of trace event types to be filtered
 8259 out. It should be possible also (as OSF suggested in their ETAP trace specifications) to select the
 8260 kernel trace event types to be traced in a system-wide fashion. These two functionalities are
 8261 called the pre-filtering of trace event types.

8262 The restriction on the actual type used for the **trace_event_set_t** type is intended to guarantee
 8263 that these objects can always be assigned, have their address taken, and be passed by value as
 8264 parameters. It is not intended that this type be a structure including pointers to other data
 8265 structures, as that could impact the portability of applications performing such operations. A
 8266 reasonable implementation could be a structure containing an array of integer types.

8267 **Filtering of Trace Event Types at Runtime**

8268 Using this API, this functionality may be built, a privileged process or a privileged thread can
 8269 get trace events from the trace stream of another process or thread, and thus specify the type of
 8270 trace events to record into a file, using methods and interfaces out of the scope of
 8271 IEEE Std 1003.1-200x. This functionality, called inline filtering of trace event types, is used for
 8272 runtime analysis of trace streams.

8273 **Post-Mortem Filtering of Trace Event Types**

8274 The word *post-mortem* is used here to indicate that some unanticipated situation occurs during
 8275 execution that does not permit a pre or inline filtering of trace events and that it is necessary to
 8276 record all trace event types, to have a chance to discover the problem afterwards. When the
 8277 program stops, all the trace events recorded previously can be analyzed in order to find the
 8278 solution. This functionality could be named the post-filtering of trace event types.

8279 **Discussions about Trace Event Type-Filtering**

8280 After long discussions with the parties involved in the process of defining the trace interface, it
 8281 seems that the sensitivity to the filtering problem is different, but everybody agrees that the level
 8282 of the overhead introduced during the tracing operation depends on the filtering method
 8283 elected. If the time that it takes the trace event to be recorded can be neglected, the overhead
 8284 introduced by the filtering process can be classified as follows:

8285 Pre-filtering System and process/thread-level overhead

8286 Inline-filtering Process/thread-level overhead

8287 Post-filtering No overhead; done offline

8288 The pre-filtering could be named *critical realtime* filtering in the sense that the filtering of trace
 8289 event type is manageable at the user level so the user can lower to a minimum the filtering
 8290 overhead at some user selected level of priority for the inline filtering, or delay the filtering to
 8291 after execution for the post-filtering. The counterpart of this solution is that the size of the trace
 8292 stream must be sufficient to record all the trace events. The advantage of the pre-filtering is that
 8293 the utilization of the trace stream is optimized.

8294 Only pre-filtering is defined by IEEE Std 1003.1-200x. However, great care must be taken in
 8295 specifying pre-filtering, so that it does not impose unacceptable overhead. Moreover, it is
 8296 necessary to isolate all the functionality relative to the pre-filtering.

8297 The result of this rationale is to define a new option, the Trace Event Filter option, not
 8298 necessarily implemented in small realtime systems, where system overhead is minimized to the
 8299 extent possible.

8300 *B.2.11.7 Tracing, pthread API*

8301 The objective to be able to control tracing for individual threads may be in conflict with the
 8302 efficiency expected in threads with a `contentionscope` attribute of
 8303 `PTHREAD_SCOPE_PROCESS`. For these threads, context switches from one thread that has
 8304 tracing enabled to another thread that has tracing disabled may require a kernel call to inform
 8305 the kernel whether it has to trace system events executed by that thread or not. For this reason, it
 8306 was proposed that the ability to enable or disable tracing for `PTHREAD_SCOPE_PROCESS`
 8307 threads be made optional, through the introduction of a Trace Scope Process option. A trace
 8308 implementation which did not implement the Trace Scope Process option would not honor the
 8309 tracing-state attribute of a thread with `PTHREAD_SCOPE_PROCESS`; it would, however, honor
 8310 the tracing-state attribute of a thread with `PTHREAD_SCOPE_SYSTEM`. This proposal was
 8311 rejected as:

- 8312 1. Removing desired functionality (per-thread trace control)
- 8313 2. Introducing counter-intuitive behavior for the tracing-state attribute
- 8314 3. Mixing logically orthogonal ideas (thread scheduling and thread tracing)
- 8315 [Objective 4]

8316 Finally, to solve this complex issue, this API does not provide `pthread_gettracingstate()`,
 8317 `pthread_settracingstate()`, `pthread_attr_gettracingstate()`, and `pthread_attr_settracingstate()`
 8318 interfaces. These interfaces force the thread implementation to add to the weight of the thread
 8319 and cause a revision of the threads libraries, just to support tracing. Worse yet,
 8320 `posix_trace_userevent()` must always test this per-thread variable even in the common case where
 8321 it is not used at all. Per-thread tracing is easy to implement using existing interfaces where
 8322 necessary; see the following example.

8323 **Example**

```
8324 /* Caution. Error checks omitted */
8325 static pthread_key_t my_key;
8326 static trace_event_id_t my_event_id;
8327 static pthread_once_t my_once = PTHREAD_ONCE_INIT;

8328 void my_init(void)
8329 {
8330     (void) pthread_key_create(&my_key, NULL);
8331     (void) posix_trace_eventid_open("my", &my_event_id);
8332 }

8333 int get_trace_flag(void)
8334 {
8335     pthread_once(&my_once, my_init);
8336     return (pthread_getspecific(my_key) != NULL);
8337 }

8338 void set_trace_flag(int f)
8339 {
8340     pthread_once(&my_once, my_init);
8341     pthread_setspecific(my_key, f? &my_event_id: NULL);
8342 }

8343 fn()
8344 {
8345     if (get_trace_flag())
```

```
8346         posix_trace_event(my_event_id, ...)
8347     }
```

8348 The above example does not implement third-party state setting, but it is also implementable
8349 with some more work, yet the extra functionality is rarely needed.

8350 Lastly, per-thread tracing works poorly for threads with PTHREAD_SCOPE_PROCESS
8351 contention scope. These “library” threads have minimal interaction with the kernel and would
8352 have to explicitly set the attributes whenever they are context switched to a new kernel thread in
8353 order to trace system events. Such state was explicitly avoided in POSIX threads to keep
8354 PTHREAD_SCOPE_PROCESS threads lightweight.

8355 The reason that keeping PTHREAD_SCOPE_PROCESS threads lightweight is important is that
8356 such threads can be used not just for simple multi-processors but also for coroutine style
8357 programming (such as discrete event simulation) without inventing a new threads paradigm.
8358 Adding extra runtime cost to thread context switches will make using POSIX threads less
8359 attractive in these situations.

8360 *B.2.11.8 Rationale on Triggering*

8361 The ability to start or stop tracing based on the occurrence of specific trace event types has been
8362 proposed as a parallel to similar functionality appearing in logic analyzers. Such triggering, in
8363 order to be very useful, should be based not only on the trace event type, but on trace event-
8364 specific data, including tests of user-specified fields for matching or threshold values.

8365 Such a facility is unnecessary where the buffering of the stream is not a constraint, since such
8366 checks can be performed offline during post-mortem analysis.

8367 For example, a large system could incorporate a daemon utility to collect the trace records from
8368 memory buffers and spool them to secondary storage for later analysis. In the instances where
8369 resources are truly limited, such as embedded applications, the application incorporation of
8370 application code to test the circumstances of a trace event and call the trace point only if needed
8371 is usually straightforward.

8372 For performance reasons, the *posix_trace_event()* function should be implemented using a macro,
8373 so if the trace is inactive, the trace event point calls are latent code and must cost no more than a
8374 scalar test.

8375 The API proposed in IEEE Std 1003.1-200x does not include any triggering functionality.

8376 *B.2.11.9 Rationale on Timestamp Clock*

8377 It has been suggested that the tracing mechanism should include the possibility of specifying the
8378 clock to be used in timestamping the trace events. When application trace events must be
8379 correlated to remote trace events, such a facility could provide a global time reference not
8380 available from a local clock. Further, the application may be driven by timers based on a clock
8381 different from that used for the timestamp, and the correlation of the trace to those untraced
8382 timer activities could be an important part of the analysis of the application.

8383 However, the tracing mechanism needs to be fast and just the provision of such an option can
8384 materially affect its performance. Leaving aside the performance costs of reading some clocks,
8385 this notion is also ill-defined when kernel trace events are to be traced by two applications
8386 making use of different tracing clocks. This can even happen within a single application where
8387 different parts of the application are served by different clocks. Another complication can occur
8388 when a clock is maintained strictly at the user level and is unavailable at the kernel level.

8389 It is felt that the benefits of a selectable trace clock do not match its costs. Applications that wish
8390 to correlate clocks other than the default tracing clock can include trace events with sample

8391 values of those other clocks, allowing correlation of timestamps from the various independent
8392 clocks. In any case, such a technique would be required when applications are sensitive to
8393 multiple clocks.

8394 *B.2.11.10 Rationale on Different Overrun Conditions*

8395 The analysis of the dynamic behavior of the trace mechanism shows that different overrun
8396 conditions may occur. The API must provide a means to manage such conditions in a portable
8397 way.

8398 **Overrun in Trace Streams Initialized with POSIX_TRACE_LOOP Policy**

8399 In this case, the user of the trace mechanism is interested in using the trace stream with
8400 POSIX_TRACE_LOOP policy to record trace events continuously, but ideally without losing any
8401 trace events. The online analyzer process must get the trace events at a mean speed equivalent to
8402 the recording speed. Should the trace stream become full, a trace stream overrun occurs. This
8403 condition is detected by getting the status of the active trace stream (function *posix_trace_get_status()*)
8404 and looking at the member *posix_stream_overrun_status* of the read **posix_stream_status**
8405 structure. In addition, two predefined trace event types are defined:

- 8406 1. The beginning of a trace overflow, to locate the beginning of an overflow when reading a
8407 trace stream
- 8408 2. The end of a trace overflow, to locate the end of an overflow, when reading a trace stream

8409 As a timestamp is associated with these predefined trace events, it is possible to know the
8410 duration of the overflow.

8411 **Overrun in Dumping Trace Streams into Trace Logs**

8412 The user lets the trace mechanism dump the trace stream initialized with
8413 POSIX_TRACE_FLUSH policy automatically into a trace log. If the dump operation is slower
8414 than the recording of trace events, the trace stream can overrun. This condition is detected by
8415 getting the status of the active trace stream (function *posix_trace_get_status()*) and looking at the
8416 member *posix_log_overrun_status* of the read **posix_stream_status** structure. This overrun
8417 indicates that the trace mechanism is not able to operate in this mode at this speed. It is the
8418 responsibility of the user to modify one of the trace parameters (the stream size or the trace
8419 event type filter, for instance) to avoid such overrun conditions, if overruns are to be prevented.
8420 The same already predefined trace event types (see **Overrun in Trace Streams Initialized with**
8421 **POSIX_TRACE_LOOP Policy**) are used to detect and to know the duration of an overflow.

8422 **Reading an Active Trace Stream**

8423 Although this trace API allows one to read an active trace stream with log while it is tracing, this
8424 feature can lead to false overflow origin interpretation: the trace log or the reader of the trace
8425 stream. Reading from an active trace stream with log is thus non-portable, and has been left
8426 unspecified.

8427 **B.2.12 Data Types**

8428 The requirement that additional types defined in this section end in “_t” was prompted by the
 8429 problem of name space pollution. It is difficult to define a type (where that type is not one
 8430 defined by IEEE Std 1003.1-200x) in one header file and use it in another without adding symbols
 8431 to the name space of the program. To allow implementors to provide their own types, all
 8432 conforming applications are required to avoid symbols ending in “_t”, which permits the
 8433 implementor to provide additional types. Because a major use of types is in the definition of
 8434 structure members, which can (and in many cases must) be added to the structures defined in
 8435 IEEE Std 1003.1-200x, the need for additional types is compelling.

8436 The types, such as **ushort** and **ulong**, which are in common usage, are not defined in
 8437 IEEE Std 1003.1-200x (although **ushort_t** would be permitted as an extension). They can be
 8438 added to `<sys/types.h>` using a feature test macro (see Section B.2.2.1 (on page 3375)). A
 8439 suggested symbol for these is `_SYSIII`. Similarly, the types like **u_short** would probably be best
 8440 controlled by `_BSD`.

8441 Some of these symbols may appear in other headers; see Section B.2.2.2 (on page 3376).

8442 **dev_t** This type may be made large enough to accommodate host-locality considerations
 8443 of networked systems.

8444 This type must be arithmetic. Earlier proposals allowed this to be non-arithmetic
 8445 (such as a structure) and provided a `samefile()` function for comparison.

8446 **gid_t** Some implementations had separated **gid_t** from **uid_t** before POSIX.1 was
 8447 completed. It would be difficult for them to coalesce them when it was
 8448 unnecessary. Additionally, it is quite possible that user IDs might be different from
 8449 group IDs because the user ID might wish to span a heterogeneous network,
 8450 where the group ID might not.

8451 For current implementations, the cost of having a separate **gid_t** will be only
 8452 lexical.

8453 **mode_t** This type was chosen so that implementations could choose the appropriate
 8454 integral type, and for compatibility with the ISO C standard. 4.3 BSD uses
 8455 **unsigned short** and the SVID uses **ushort**, which is the same. Historically, only the
 8456 low-order sixteen bits are significant.

8457 **nlink_t** This type was introduced in place of **short** for `st_nlink` (see the `<sys/stat.h>` header)
 8458 in response to an objection that **short** was too small.

8459 **off_t** This type is used only in `lseek()`, `fcntl()`, and `<sys/stat.h>`. Many implementations
 8460 would have difficulties if it were defined as anything other than **long**. Requiring
 8461 an integral type limits the capabilities of `lseek()` to four gigabytes. The ISO C
 8462 standard supplies routines that use larger types; see `fgetpos()` and `fsetpos()`. XSI-
 8463 conformant systems provide the `fseeko()` and `lseeko()` functions that use larger
 8464 types.

8465 **pid_t** The inclusion of this symbol was controversial because it is tied to the issue of the
 8466 representation of a process ID as a number. From the point of view of a
 8467 conforming application, process IDs should be “magic cookies”¹ that are produced

8468 _____

8469 1. An historical term meaning: “An opaque object, or token, of determinate size, whose significance is known only to the entity
 8470 which created it. An entity receiving such a token from the generating entity may only make such use of the ‘cookie’ as is defined
 8471 and permitted by the supplying entity.”

8472 by calls such as *fork()*, used by calls such as *waitpid()* or *kill()*, and not otherwise
8473 analyzed (except that the sign is used as a flag for certain operations).

8474 The concept of a {PID_MAX} value interacted with this in early proposals. Treating
8475 process IDs as an opaque type both removes the requirement for {PID_MAX} and
8476 allows systems to be more flexible in providing process IDs that span a large range
8477 of values, or a small one.

8478 Since the values in **uid_t**, **gid_t**, and **pid_t** will be numbers generally, and
8479 potentially both large in magnitude and sparse, applications that are based on
8480 arrays of objects of this type are unlikely to be fully portable in any case. Solutions
8481 that treat them as magic cookies will be portable.

8482 {CHILD_MAX} precludes the possibility of a “toy implementation”, where there
8483 would only be one process.

8484 **ssize_t** This is intended to be a signed analog of **size_t**. The wording is such that an
8485 implementation may either choose to use a longer type or simply to use the signed
8486 version of the type that underlies **size_t**. All functions that return **ssize_t** (*read()*
8487 and *write()*) describe as “implementation-defined” the result of an input exceeding
8488 {SSIZE_MAX}. It is recognized that some implementations might have **ints** that
8489 are smaller than **size_t**. A conforming application would be constrained not to
8490 perform I/O in pieces larger than {SSIZE_MAX}, but a conforming application
8491 using extensions would be able to use the full range if the implementation
8492 provided an extended range, while still having a single type-compatible interface.

8493 The symbols **size_t** and **ssize_t** are also required in **<unistd.h>** to minimize the
8494 changes needed for calls to *read()* and *write()*. Implementors are reminded that it
8495 must be possible to include both **<sys/types.h>** and **<unistd.h>** in the same
8496 program (in either order) without error.

8497 **uid_t** Before the addition of this type, the data types used to represent these values
8498 varied throughout early proposals. The **<sys/stat.h>** header defined these values as
8499 type **short**, the **<passwd.h>** file (now **<pwd.h>** and **<grp.h>**) used an **int**, and
8500 *getuid()* returned an **int**. In response to a strong objection to the inconsistent
8501 definitions, all the types to were switched to **uid_t**.

8502 In practice, those historical implementations that use varying types of this sort can
8503 typedef **uid_t** to **short** with no serious consequences.

8504 The problem associated with this change concerns object compatibility after
8505 structure size changes. Since most implementations will define **uid_t** as a short, the
8506 only substantive change will be a reduction in the size of the **passwd** structure.
8507 Consequently, implementations with an overriding concern for object
8508 compatibility can pad the structure back to its current size. For that reason, this
8509 problem was not considered critical enough to warrant the addition of a separate
8510 type to POSIX.1.

8511 The types **uid_t** and **gid_t** are magic cookies. There is no {UID_MAX} defined by
8512 POSIX.1, and no structure imposed on **uid_t** and **gid_t** other than that they be
8513 positive arithmetic types. (In fact, they could be **unsigned char**.) There is no
8514 maximum or minimum specified for the number of distinct user or group IDs.

8515 **B.3 System Interfaces**

8516 See the RATIONALE sections on the individual reference pages.

8517 **B.3.1 Examples for Spawn**8518 The following long examples are provided in the Rationale (Informative) volume of
8519 IEEE Std 1003.1-200x as a supplement to the reference page for *spawn()*.8520 **Example Library Implementation of Spawn**8521 The *posix_spawn()* or *posix_spawnnp()* functions provide the following:

- 8522 • Simply start a process executing a process image. This is the simplest application for process
8523 creation, and it may cover most executions of POSIX *fork()*.
- 8524 • Support I/O redirection, including pipes.
- 8525 • Run the child under a user and group ID in the domain of the parent.
- 8526 • Run the child at any priority in the domain of the parent.

8527 The *posix_spawn()* or *posix_spawnnp()* functions do not cover every possible use of the *fork()*
8528 function, but they do span the common applications: typical use by a shell and a login utility.8529 The price for an application is that before it calls *posix_spawn()* or *posix_spawnnp()*, the parent
8530 must adjust to a state that *posix_spawn()* or *posix_spawnnp()* can map to the desired state for the
8531 child. Environment changes require the parent to save some of its state and restore it afterwards.
8532 The effective behavior of a successful invocation of *posix_spawn()* is as if the operation were
8533 implemented with POSIX operations as follows:

```

8534 #include <sys/types.h>
8535 #include <stdlib.h>
8536 #include <stdio.h>
8537 #include <unistd.h>
8538 #include <sched.h>
8539 #include <fcntl.h>
8540 #include <signal.h>
8541 #include <errno.h>
8542 #include <string.h>
8543 #include <signal.h>

8544 /* #include <spawn.h>*/
8545 /*****
8546 /* Things that could be defined in spawn.h */
8547 /*****
8548 typedef struct
8549     {
8550     short posix_attr_flags;
8551     #define POSIX_SPAWN_SETPGROUP          0x1
8552     #define POSIX_SPAWN_SETSIGMASK       0x2
8553     #define POSIX_SPAWN_SETSIGDEF       0x4
8554     #define POSIX_SPAWN_SETSCHEDULER    0x8
8555     #define POSIX_SPAWN_SETSCHEDPARAM  0x10
8556     #define POSIX_SPAWN_RESETIDS        0x20
8557     pid_t posix_attr_pgroup;
8558     sigset_t posix_attr_sigmask;
8559     sigset_t posix_attr_sigdefault;

```

```

8560         int posix_attr_schedpolicy;
8561         struct sched_param posix_attr_schedparam;
8562     } posix_spawnattr_t;

8563     typedef char *posix_spawn_file_actions_t;

8564     int posix_spawn_file_actions_init(
8565         posix_spawn_file_actions_t *file_actions);
8566     int posix_spawn_file_actions_destroy(
8567         posix_spawn_file_actions_t *file_actions);
8568     int posix_spawn_file_actions_addclose(
8569         posix_spawn_file_actions_t *file_actions, int fildes);
8570     int posix_spawn_file_actions_adddup2(
8571         posix_spawn_file_actions_t *file_actions, int fildes,
8572         int newfildes);
8573     int posix_spawn_file_actions_addopen(
8574         posix_spawn_file_actions_t *file_actions, int fildes,
8575         const char *path, int oflag, mode_t mode);
8576     int posix_spawnattr_init(posix_spawnattr_t *attr);
8577     int posix_spawnattr_destroy(posix_spawnattr_t *attr);
8578     int posix_spawnattr_getflags(const posix_spawnattr_t *attr, short *lags);
8579     int posix_spawnattr_setflags(posix_spawnattr_t *attr, short flags);
8580     int posix_spawnattr_getpgroup(const posix_spawnattr_t *attr,
8581         pid_t *pgroup);
8582     int posix_spawnattr_setpgroup(posix_spawnattr_t *attr, pid_t pgroup);
8583     int posix_spawnattr_getschedpolicy(const posix_spawnattr_t *attr,
8584         int *schedpolicy);
8585     int posix_spawnattr_setschedpolicy(posix_spawnattr_t *attr,
8586         int schedpolicy);
8587     int posix_spawnattr_getschedparam(const posix_spawnattr_t *attr,
8588         struct sched_param *schedparam);
8589     int posix_spawnattr_setschedparam(posix_spawnattr_t *attr,
8590         const struct sched_param *schedparam);
8591     int posix_spawnattr_getsigmask(const posix_spawnattr_t *attr,
8592         sigset_t *sigmask);
8593     int posix_spawnattr_setsigmask(posix_spawnattr_t *attr,
8594         const sigset_t *sigmask);
8595     int posix_spawnattr_getdefault(const posix_spawnattr_t *attr,
8596         sigset_t *sigdefault);
8597     int posix_spawnattr_setsigdefault(posix_spawnattr_t *attr,
8598         const sigset_t *sigdefault);
8599     int posix_spawn(pid_t *pid, const char *path,
8600         const posix_spawn_file_actions_t *file_actions,
8601         const posix_spawnattr_t *attrp, char * const argv[],
8602         char * const envp[]);
8603     int posix_spawnnp(pid_t *pid, const char *file,
8604         const posix_spawn_file_actions_t *file_actions,
8605         const posix_spawnattr_t *attrp, char * const argv[],
8606         char * const envp[]);

8607     /*****
8608     /* Example posix_spawn() library routine */
8609     /*****
8610     int posix_spawn(pid_t *pid,

```

```

8611     const char *path,
8612     const posix_spawn_file_actions_t *file_actions,
8613     const posix_spawnattr_t *attrp,
8614     char * const argv[],
8615     char * const envp[])
8616 {
8617     /* Create process */
8618     if((*pid=fork()) == (pid_t)0)
8619     {
8620         /* This is the child process */
8621         /* Worry about process group */
8622         if(attrp->posix_attr_flags & POSIX_SPAWN_SETPGROUP)
8623         {
8624             /* Override inherited process group */
8625             if(setpgid(0, attrp->posix_attr_pgroup) != 0)
8626             {
8627                 /* Failed */
8628                 exit(127);
8629             }
8630         }
8631
8632         /* Worry about process signal mask */
8633         if(attrp->posix_attr_flags & POSIX_SPAWN_SETSIGMASK)
8634         {
8635             /* Set the signal mask (can't fail) */
8636             sigprocmask(SIG_SETMASK , &attrp->posix_attr_sigmask,
8637                 NULL);
8638         }
8639
8640         /* Worry about resetting effective user and group IDs */
8641         if(attrp->posix_attr_flags & POSIX_SPAWN_RESETEIDS)
8642         {
8643             /* None of these can fail for this case. */
8644             setuid(getuid());
8645             setgid(getgid());
8646         }
8647
8648         /* Worry about defaulted signals */
8649         if(attrp->posix_attr_flags & POSIX_SPAWN_SETSIGDEF)
8650         {
8651             struct sigaction deflt;
8652             sigset_t all_signals;
8653
8654             int s;
8655
8656             /* Construct default signal action */
8657             deflt.sa_handler = SIG_DFL;
8658             deflt.sa_flags = 0;
8659
8660             /* Construct the set of all signals */
8661             sigfillset(&all_signals);
8662
8663             /* Loop for all signals */
8664             for(s=0; sigismember(&all_signals,s); s++)
8665             {
8666                 /* Signal to be defaulted? */

```



```

8660         if(sigismember(&attrp->posix_attr_sigdefault,s))
8661             {
8662                 /* Yes; default this signal */
8663                 if(sigaction(s, &deflt, NULL) == -1)
8664                     {
8665                         /* Failed */
8666                         exit(127);
8667                     }
8668             }
8669     }
8670 }

8671 /* Worry about the fds if we are to map them */
8672 if(file_actions != NULL)
8673     {
8674     /* Loop for all actions in object file_actions */
8675     /*(implementation dives beneath abstraction)*/
8676     char *p = *file_actions;
8677     while(*p != ' ')
8678         {
8679             if(strncmp(p,"close(",6) == 0)
8680                 {
8681                     int fd;
8682                     if(sscanf(p+6,"%d",&fd) != 1)
8683                         {
8684                             exit(127);
8685                         }
8686                     if(close(fd) == -1) exit(127);
8687                 }
8688             else if(strncmp(p,"dup2(",5) == 0)
8689                 {
8690                     int fd,newfd;
8691                     if(sscanf(p+5,"%d,%d",&fd,&newfd) != 2)
8692                         {
8693                             exit(127);
8694                         }
8695                     if(dup2(fd, newfd) == -1) exit(127);
8696                 }
8697             else if(strncmp(p,"open(",5) == 0)
8698                 {
8699                     int fd,oflag;
8700                     mode_t mode;
8701                     int tempfd;
8702                     char path[1000]; /* Should be dynamic */
8703                     char *q;
8704                     if(sscanf(p+5,"%d",&fd) != 1)
8705                         {
8706                             exit(127);
8707                         }
8708                     p = strchr(p, ' ') + 1;
8709                     q = strchr(p, '*');
8710                     if(q == NULL) exit(127);
8711                     strncpy(path, p, q-p);

```

```

8712         path[q-p] = ' ';
8713     if(sscanf(q+1,"%o,%o",&oflag,&mode)!=2)
8714     {
8715         exit(127);
8716     }
8717     if(close(fd) == -1)
8718     {
8719         if(errno != EBADF) exit(127);
8720     }
8721     tempfd = open(path, oflag, mode);
8722     if(tempfd == -1) exit(127);
8723     if(tempfd != fd)
8724     {
8725         if(dup2(tempfd,fd) == -1)
8726         {
8727             exit(127);
8728         }
8729         if(close(tempfd) == -1)
8730         {
8731             exit(127);
8732         }
8733     }
8734     }
8735     else
8736     {
8737         exit(127);
8738     }
8739     p = strchr(p, ' ') + 1;
8740 }
8741 }
8742 /* Worry about setting new scheduling policy and parameters */
8743 if(attrp->posix_attr_flags & POSIX_SPAWN_SETSCHEDULER)
8744 {
8745     if(sched_setscheduler(0, attrp->posix_attr_schedpolicy,
8746         &attrp->posix_attr_schedparam) == -1)
8747     {
8748         exit(127);
8749     }
8750 }
8751 /* Worry about setting only new scheduling parameters */
8752 if(attrp->posix_attr_flags & POSIX_SPAWN_SETSCHEDPARAM)
8753 {
8754     if(sched_setparam(0, &attrp->posix_attr_schedparam)==-1)
8755     {
8756         exit(127);
8757     }
8758 }
8759 /* Now execute the program at path */
8760 /* Any fd that still has FD_CLOEXEC set will be closed */
8761 execve(path, argv, envp);
8762 exit(127); /* exec failed */

```

```

8763     }
8764     else
8765     {
8766         /* This is the parent (calling) process */
8767         if(*pid == (pid_t)-1) return errno;
8768         return 0;
8769     }
8770 }

8771 /*****
8772  /* Here is a crude but effective implementation of the */
8773  /* file action object operators which store actions as */
8774  /* concatenated token separated strings.          */
8775  *****/
8776 /* Create object with no actions. */
8777 int posix_spawn_file_actions_init(
8778     posix_spawn_file_actions_t *file_actions)
8779 {
8780     *file_actions = malloc(sizeof(char));
8781     if(*file_actions == NULL) return ENOMEM;
8782     strcpy(*file_actions, "");
8783     return 0;
8784 }

8785 /* Free object storage and make invalid. */
8786 int posix_spawn_file_actions_destroy(
8787     posix_spawn_file_actions_t *file_actions)
8788 {
8789     free(*file_actions);
8790     *file_actions = NULL;
8791     return 0;
8792 }

8793 /* Add a new action string to object. */
8794 static int add_to_file_actions(
8795     posix_spawn_file_actions_t *file_actions,
8796     char *new_action)
8797 {
8798     *file_actions = realloc
8799         (*file_actions, strlen(*file_actions)+strlen(new_action)+1);
8800     if(*file_actions == NULL) return ENOMEM;
8801     strcat(*file_actions, new_action);
8802     return 0;
8803 }

8804 /* Add a close action to object. */
8805 int posix_spawn_file_actions_addclose(
8806     posix_spawn_file_actions_t *file_actions, int fildes)
8807 {
8808     char temp[100];
8809     sprintf(temp, "close(%d)", fildes);
8810     return add_to_file_actions(file_actions, temp);
8811 }

```

```

8812     /* Add a dup2 action to object. */
8813     int posix_spawn_file_actions_adddup2(
8814         posix_spawn_file_actions_t *file_actions, int fildes,
8815         int newfildes)
8816     {
8817         char temp[100];
8818         sprintf(temp, "dup2(%d,%d)", fildes, newfildes);
8819         return add_to_file_actions(file_actions, temp);
8820     }

8821     /* Add an open action to object. */
8822     int posix_spawn_file_actions_addopen(
8823         posix_spawn_file_actions_t *file_actions, int fildes,
8824         const char *path, int oflag, mode_t mode)
8825     {
8826         char temp[100];
8827         sprintf(temp, "open(%d,%s*%o,%o)", fildes, path, oflag, mode);
8828         return add_to_file_actions(file_actions, temp);
8829     }

8830     /******
8831     /* Here is a crude but effective implementation of the */
8832     /* spawn attributes object functions which manipulate */
8833     /* the individual attributes. */
8834     /******
8835     /* Initialize object with default values. */
8836     int posix_spawnattr_init(
8837         posix_spawnattr_t *attr)
8838     {
8839         attr->posix_attr_flags=0;
8840         attr->posix_attr_pgroup=0;
8841         /* Default value of signal mask is the parent's signal mask; */
8842         /* other values are also allowed */
8843         sigprocmask(0,NULL,&attr->posix_attr_sigmask);
8844         sigemptyset(&attr->posix_attr_sigdefault);
8845         /* Default values of scheduling attr inherited from the parent; */
8846         /* other values are also allowed */
8847         attr->posix_attr_schedpolicy=sched_getscheduler(0);
8848         sched_getparam(0,&attr->posix_attr_schedparam);
8849         return 0;
8850     }

8851     int posix_spawnattr_destroy(posix_spawnattr_t *attr)
8852     {
8853         /* No action needed */
8854         return 0;
8855     }

8856     int posix_spawnattr_getflags(const posix_spawnattr_t *attr,
8857         short *flags)
8858     {
8859         *flags=attr->posix_attr_flags;
8860         return 0;
8861     }

```

```
8862     int posix_spawnattr_setflags(posix_spawnattr_t *attr, short flags)
8863     {
8864         attr->posix_attr_flags=flags;
8865         return 0;
8866     }

8867     int posix_spawnattr_getpgroup(const posix_spawnattr_t *attr,
8868         pid_t *pgroup)
8869     {
8870         *pgroup=attr->posix_attr_pgroup;
8871         return 0;
8872     }

8873     int posix_spawnattr_setpgroup(posix_spawnattr_t *attr, pid_t pgroup)
8874     {
8875         attr->posix_attr_pgroup=pgroup;
8876         return 0;
8877     }

8878     int posix_spawnattr_getschedpolicy(const posix_spawnattr_t *attr,
8879         int *schedpolicy)
8880     {
8881         *schedpolicy=attr->posix_attr_schedpolicy;
8882         return 0;
8883     }

8884     int posix_spawnattr_setschedpolicy(posix_spawnattr_t *attr,
8885         int schedpolicy)
8886     {
8887         attr->posix_attr_schedpolicy=schedpolicy;
8888         return 0;
8889     }

8890     int posix_spawnattr_getschedparam(const posix_spawnattr_t *attr,
8891         struct sched_param *schedparam)
8892     {
8893         *schedparam=attr->posix_attr_schedparam;
8894         return 0;
8895     }

8896     int posix_spawnattr_setschedparam(posix_spawnattr_t *attr,
8897         const struct sched_param *schedparam)
8898     {
8899         attr->posix_attr_schedparam=*schedparam;
8900         return 0;
8901     }

8902     int posix_spawnattr_getsigmask(const posix_spawnattr_t *attr,
8903         sigset_t *sigmask)
8904     {
8905         *sigmask=attr->posix_attr_sigmask;
8906         return 0;
8907     }

8908     int posix_spawnattr_setsigmask(posix_spawnattr_t *attr,
8909         const sigset_t *sigmask)
```

```

8910     {
8911     attr->posix_attr_sigmask=*sigmask;
8912     return 0;
8913     }

8914     int posix_spawnattr_getsigdefault(const posix_spawnattr_t *attr,
8915     sigset_t *sigdefault)
8916     {
8917     *sigdefault=attr->posix_attr_sigdefault;
8918     return 0;
8919     }

8920     int posix_spawnattr_setsigdefault(posix_spawnattr_t *attr,
8921     const sigset_t *sigdefault)
8922     {
8923     attr->posix_attr_sigdefault=*sigdefault;
8924     return 0;
8925     }

```

8926 **I/O Redirection with Spawn**

8927 I/O redirection with *posix_spawn()* or *posix_spawnp()* is accomplished by crafting a *file_actions*
8928 argument to effect the desired redirection. Such a redirection follows the general outline of the
8929 following example:

```

8930     /* To redirect new standard output (fd 1) to a file, */
8931     /* and redirect new standard input (fd 0) from my fd socket_pair[1], */
8932     /* and close my fd socket_pair[0] in the new process. */
8933     posix_spawn_file_actions_t file_actions;
8934     posix_spawn_file_actions_init(&file_actions);
8935     posix_spawn_file_actions_addopen(&file_actions, 1, "newout", ...);
8936     posix_spawn_file_actions_dup2(&file_actions, socket_pair[1], 0);
8937     posix_spawn_file_actions_close(&file_actions, socket_pair[0]);
8938     posix_spawn_file_actions_close(&file_actions, socket_pair[1]);
8939     posix_spawn(..., &file_actions, ...);
8940     posix_spawn_file_actions_destroy(&file_actions);

```

8941 **Spawning a Process Under a New User ID**

8942 Spawning a process under a new user ID follows the outline shown in the following example:

```

8943     Save = getuid();
8944     setuid(newid);
8945     posix_spawn(...);
8946     setuid(Save);

```


8948 / *Rationale (Informative)*

8949 **Part C:**

8950 **Shell and Utilities**

8951 *The Open Group*
8952 *The Institute of Electrical and Electronics Engineers, Inc.*

Rationale for Shell and Utilities

8953

8954 **C.1 Introduction**8955 **C.1.1 Scope**

8956 Refer to Section A.1.1 (on page 3293).

8957 **C.1.2 Conformance**

8958 Refer to Section A.2 (on page 3299).

8959 **C.1.3 Normative References**

8960 There is no additional rationale provided for this section. |

8961 **C.1.4 Change History** |8962 The change history is provided as an informative section, to track changes from previous issues |
8963 of IEEE Std 1003.1-200x. |8964 The following sections describe changes made to the Shell and Utilities volume of |
8965 IEEE Std 1003.1-200x since Issue 5 of the base document. The CHANGE HISTORY section for |
8966 each utility describes technical changes made to that utility since Issue 5. Changes between |
8967 earlier issues of the base document and Issue 5 are not included. |8968 The change history between Issue 5 and Issue 6 also lists the changes since the |
8969 ISO POSIX-2: 1993 standard. |8970 **Changes from Issue 5 to Issue 6 (IEEE Std 1003.1-200x)** |8971 The following list summarizes the major changes that were made in the Shell and Utilities |
8972 volume of IEEE Std 1003.1-200x from Issue 5 to Issue 6: |8973 • This volume of IEEE Std 1003.1-200x is extensively revised so it can be both an IEEE POSIX |
8974 Standard and an Open Group Technical Standard. |

8975 • The terminology has been reworked to meet the style requirements. |

8976 • Shading notation and margin codes are introduced for identification of options within the |
8977 volume. |8978 • This volume of IEEE Std 1003.1-200x is updated to mandate support of FIPS 151-2. The |
8979 following changes were made: |8980 — Support is mandated for the capabilities associated with the following symbolic |
8981 constants: |8982 `_POSIX_CHOWN_RESTRICTED` |8983 `_POSIX_JOB_CONTROL` |8984 `_POSIX_SAVED_IDS` |8985 — In the environment for the login shell, the environment variables *LOGNAME* and *HOME* |
8986 shall be defined and have the properties described in the Base Definitions volume of |

- 8987 IEEE Std 1003.1-200x, Chapter 7, Locale.
- 8988 • This volume of IEEE Std 1003.1-200x is updated to align with some features of the Single
 - 8989 UNIX Specification.
 - 8990 • A new section on Utility Limits is added.
 - 8991 • A section on the Relationships to Other Documents is added.
 - 8992 • Concepts and definitions have been moved to a separate volume.
 - 8993 • A RATIONALE section is added to each reference page.
 - 8994 • The *c99* utility is added as a replacement for *c89*, which is withdrawn in this issue.
 - 8995 • IEEE Std 1003.2d-1994 is incorporated, adding the *qalter*, *qdel*, *qhold*, *qmove*, *qmsg*, *qrerun*, *qrls*,
 - 8996 *qselect*, *qsig*, *qstat*, and *qsub* utilities.
 - 8997 • IEEE P1003.2b draft standard is incorporated, making extensive updates and adding the *iconv*
 - 8998 utility.
 - 8999 • IEEE PASC Interpretations are applied.
 - 9000 • The Open Groups corrigenda and resolutions are applied.

9001 **New Features in Issue 6**

9002 The following table lists the new utilities introduced since the ISO POSIX-2:1993 standard (as
9003 modified by IEEE Std 1003.2d-1994). These are all part of the XSI extension.

9004

9005

9006

9007

9008

9009

9010

9011

New Utilities in Issue 6				
<i>admin</i>	<i>fuser</i>	<i>link</i>	<i>tsort</i>	<i>uustat</i>
<i>cal</i>	<i>genccat</i>	<i>m4</i>	<i>ulimit</i>	<i>uux</i>
<i>cflow</i>	<i>get</i>	<i>nl</i>	<i>uncompress</i>	<i>val</i>
<i>compress</i>	<i>hash</i>	<i>pr</i>	<i>unget</i>	<i>what</i>
<i>cxref</i>	<i>ipcrm</i>	<i>sact</i>	<i>unlink</i>	<i>zcat</i>
<i>delta</i>	<i>ipcs</i>	<i>sccs</i>	<i>uucp</i>	

9012 **C.1.5 Terminology**

9013 Refer to Section A.1.4 (on page 3295).

9014 **C.1.6 Definitions**

9015 Refer to Section A.3 (on page 3302).

9016 **C.1.7 Relationship to Other Documents**

9017 *C.1.7.1 System Interfaces*

9018 It has been pointed out that the Shell and Utilities volume of IEEE Std 1003.1-200x assumes that
9019 a great deal of functionality from the System Interfaces volume of IEEE Std 1003.1-200x is
9020 present, but never states exactly how much (and strictly does not need to since both are
9021 mandated on a conforming system). This section is an attempt to clarify the assumptions.

9022 *C.1.7.2 Concepts Derived from the ISO C Standard*

9023 This section was introduced to address the issue that there was insufficient detail presented by
 9024 such utilities as *awk* or *sh* about their procedural control statements and their methods of
 9025 performing arithmetic functions.

9026 The ISO C standard was selected as a model because most historical implementations of the
 9027 standard utilities were written in C. Thus, it was more likely that they would act in the desired
 9028 manner without modification.

9029 Using the ISO C standard is primarily a notational convenience so that the many procedural
 9030 languages in the Shell and Utilities volume of IEEE Std 1003.1-200x would not have to be
 9031 rigorously described in every aspect. Its selection does not require that the standard utilities be
 9032 written in Standard C; they could be written in Common Usage C, Ada, Pascal, assembler
 9033 language, or anything else.

9034 The sizes of the various numeric values refer to C-language data types that are allowed to be
 9035 different sizes by the ISO C standard. Thus, like a C-language application, a shell application
 9036 cannot rely on their exact size. However, it can rely on their minimum sizes expressed in the
 9037 ISO C standard, such as {LONG_MAX} for a **long** type.

9038 **C.1.8 Portability**

9039 Refer to Section A.1.5 (on page 3298).

9040 *C.1.8.1 Codes*

9041 Refer to Section A.1.5.1 (on page 3298).

9042 **C.1.9 Utility Limits**

9043 This section grew out of an idea that originated with the original POSIX.1, in the tables of system
 9044 limits for the *sysconf()* and *pathconf()* functions. The idea being that a conforming application
 9045 can be written to use the most restrictive values that a minimal system can provide, but it should
 9046 not have to. The values provided represent compromises so that some vendors can use
 9047 historically limited versions of UNIX system utilities. They are the highest values that a strictly
 9048 conforming application can assume, given no other information.

9049 However, by using the *getconf* utility or the *sysconf()* function, the elegant application can be
 9050 tailored to more liberal values on some of the specific instances of specific implementations.

9051 There is no explicitly stated requirement that an implementation provide finite limits for any of
 9052 these numeric values; the implementation is free to provide essentially unbounded capabilities
 9053 (where it makes sense), stopping only at reasonable points such as {ULONG_MAX} (from the
 9054 ISO C standard). Therefore, applications desiring to tailor themselves to the values on a
 9055 particular implementation need to be ready for possibly huge values; it may not be a good idea
 9056 to allocate blindly a buffer for an input line based on the value of {LINE_MAX}, for instance.
 9057 However, unlike the System Interfaces volume of IEEE Std 1003.1-200x, there is no set of limits
 9058 that return a special indication meaning “unbounded”. The implementation should always
 9059 return an actual number, even if the number is very large.

9060 The statement:

9061 “It is not guaranteed that the application ...”

9062 is an indication that many of these limits are designed to ensure that implementors design their
 9063 utilities without arbitrary constraints related to unimaginative programming. There are certainly
 9064 conditions under which combinations of options can cause failures that would not render an

9065 implementation non-conforming. For example, {EXPR_NEST_MAX} and {ARG_MAX} could
 9066 collide when expressions are large; combinations of {BC_SCALE_MAX} and {BC_DIM_MAX}
 9067 could exceed virtual memory.

9068 In the Shell and Utilities volume of IEEE Std 1003.1-200x, the notion of a limit being guaranteed
 9069 for the process lifetime, as it is in the System Interfaces volume of IEEE Std 1003.1-200x, is not as
 9070 useful to a shell script. The *getconf* utility is probably a process itself, so the guarantee would be
 9071 without value. Therefore, the Shell and Utilities volume of IEEE Std 1003.1-200x requires the
 9072 guarantee to be for the session lifetime. This will mean that many vendors will either return very
 9073 conservative values or possibly implement *getconf* as a built-in.

9074 It may seem confusing to have limits that apply only to a single utility grouped into one global
 9075 section. However, the alternative, which would be to disperse them out into their utility
 9076 description sections, would cause great difficulty when *sysconf()* and *getconf* were described.
 9077 Therefore, the standard developers chose the global approach.

9078 Each language binding could provide symbol names that are slightly different from those shown |
 9079 here. For example, the C-Language Binding option adds a leading underscore to the symbols as a |
 9080 prefix.

9081 The following comments describe selection criteria for the symbols and their values:

9082 {ARG_MAX}

9083 This is defined by the System Interfaces volume of IEEE Std 1003.1-200x. Unfortunately, it is
 9084 very difficult for a conforming application to deal with this value, as it does not know how |
 9085 much of its argument space is being consumed by the environment variables of the user.

9086 {BC_BASE_MAX}

9087 {BC_DIM_MAX}

9088 {BC_SCALE_MAX}

9089 These were originally one value, {BC_SCALE_MAX}, but it was unreasonable to link all
 9090 three concepts into one limit.

9091 {CHILD_MAX}

9092 This is defined by the System Interfaces volume of IEEE Std 1003.1-200x.

9093 {COLL_WEIGHTS_MAX}

9094 The weights assigned to **order** can be considered as “passes” through the collation
 9095 algorithm.

9096 {EXPR_NEST_MAX}

9097 The value for expression nesting was borrowed from the ISO C standard.

9098 {LINE_MAX}

9099 This is a global limit that affects all utilities, unless otherwise noted. The {MAX_CANON}
 9100 value from the System Interfaces volume of IEEE Std 1003.1-200x may further limit input
 9101 lines from terminals. The {LINE_MAX} value was the subject of much debate and is a
 9102 compromise between those who wished to have unlimited lines and those who understood
 9103 that many historical utilities were written with fixed buffers. Frequently, utility writers
 9104 selected the UNIX system constant BUFSIZ to allocate these buffers; therefore, some utilities
 9105 were limited to 512 bytes for I/O lines, while others achieved 4 096 bytes or greater.

9106 It should be noted that {LINE_MAX} applies only to input line length; there is no
 9107 requirement in IEEE Std 1003.1-200x that limits the length of output lines. Utilities such as
 9108 *awk*, *sed*, and *paste* could theoretically construct lines longer than any of the input lines they
 9109 received, depending on the options used or the instructions from the application. They are
 9110 not required to truncate their output to {LINE_MAX}. It is the responsibility of the
 9111 application to deal with this. If the output of one of those utilities is to be piped into another

9112 of the standard utilities, line length restrictions will have to be considered; the *fold* utility,
 9113 among others, could be used to ensure that only reasonable line lengths reach utilities or
 9114 applications.

9115 {LINK_MAX}

9116 This is defined by the System Interfaces volume of IEEE Std 1003.1-200x.

9117 {MAX_CANON}

9118 {MAX_INPUT}

9119 {NAME_MAX}

9120 {NGROUPS_MAX}

9121 {OPEN_MAX}

9122 {PATH_MAX}

9123 {PIPE_BUF}

9124 These limits are defined by the System Interfaces volume of IEEE Std 1003.1-200x. Note that
 9125 the byte lengths described by some of these values continue to represent bytes, even if the
 9126 applicable character set uses a multi-byte encoding.

9127 {RE_DUP_MAX}

9128 The value selected is consistent with historical practice. Although the name implies that it
 9129 applies to all REs, only BREs use the interval notation $\{m,n\}$ addressed by this limit.

9130 {POSIX2_SYMLINKS}

9131 The {POSIX2_SYMLINKS} variable indicates that the underlying operating system supports
 9132 the creation of symbolic links in specific directories. Many of the utilities defined in
 9133 IEEE Std 1003.1-200x that deal with symbolic links do not depend on this value. For
 9134 example, a utility that follows symbolic links (or does not, as the case may be) will only be
 9135 affected by a symbolic link if it encounters one. Presumably, a file system that does not
 9136 support symbolic links will not contain any. This variable does affect such utilities as *ln -s*
 9137 and *pax* that attempt to create symbolic links.

9138 {POSIX2_SYMLINKS} was developed even though there is no comparable configuration
 9139 value for the system interfaces.

9140 There are different limits associated with command lines and input to utilities, depending on the
 9141 method of invocation. In the case of a C program *exec-ing* a utility, {ARG_MAX} is the
 9142 underlying limit. In the case of the shell reading a script and *exec-ing* a utility, {LINE_MAX}
 9143 limits the length of lines the shell is required to process, and {ARG_MAX} will still be a limit. If a
 9144 user is entering a command on a terminal to the shell, requesting that it invoke the utility,
 9145 {MAX_INPUT} may restrict the length of the line that can be given to the shell to a value below
 9146 {LINE_MAX}.

9147 When an option is supported, *getconf* returns a value of 1. For example, when C development is
 9148 supported:

```
9149     if [ "$(getconf POSIX2_C_DEV)" -eq 1 ]; then
9150         echo C supported
9151     fi
```

9152 The *sysconf()* function in the C-Language Binding option would return 1.

9153 The following comments describe selection criteria for the symbols and their values:

9154 POSIX2_C_BIND

9155 POSIX2_C_DEV

9156 POSIX2_FORT_DEV

9157 POSIX2_FORT_RUN

9158 POSIX2_SW_DEV

9159 POSIX2_UPE

9160 It is possible for some (usually privileged) operations to remove utilities that support these
 9161 options or otherwise to render these options unsupported. The header files, the *sysconf()*
 9162 function, or the *getconf* utility will not necessarily detect such actions, in which case they
 9163 should not be considered as rendering the implementation non-conforming. A test suite
 9164 should not attempt tests such as:

```
9165     rm /usr/bin/c99
9166     getconf POSIX2_C_DEV
```

9167 POSIX2_LOCALEDEF

9168 This symbol was introduced to allow implementations to restrict supported locales to only
 9169 those supplied by the implementation.

9170 C.1.10 Grammar Conventions

9171 There is no additional rationale provided for this section.

9172 C.1.11 Utility Description Defaults

9173 This section is arranged with headings in the same order as all the utility descriptions. It is a
 9174 collection of related and unrelated information concerning

- 9175 1. The default actions of utilities
- 9176 2. The meanings of notations used in IEEE Std 1003.1-200x that are specific to individual
 9177 utility sections

9178 Although this material may seem out of place here, it is important that this information appear
 9179 before any of the utilities to be described later.

9180 NAME

9181 There is no additional rationale provided for this section.

9182 SYNOPSIS

9183 There is no additional rationale provided for this section.

9184 DESCRIPTION

9185 There is no additional rationale provided for this section.

9186 OPTIONS

9187 Although it has not always been possible, the standard developers tried to avoid repeating
 9188 information to reduce the risk that duplicate explanations could each be modified differently.

9189 The need to recognize `--` is required because conforming applications need to shield their |
 9190 operands from any arbitrary options that the implementation may provide as an extension. For |
 9191 example, if the standard utility *foo* is listed as taking no options, and the application needed to
 9192 give it a pathname with a leading hyphen, it could safely do it as:

```
9193     foo -- -myfile
```

9194 and avoid any problems with `-m` used as an extension.

9195 **OPERANDS**

9196 The usage of `-` is never shown in the SYNOPSIS. Similarly, the usage of `--` is never shown.

9197 The requirement for processing operands in command-line order is to avoid a “WeirdNIX”
 9198 utility that might choose to sort the input files alphabetically, by size, or by directory order.
 9199 Although this might be acceptable for some utilities, in general the programmer has a right to
 9200 know exactly what order will be chosen.

9201 Some of the standard utilities take multiple *file* operands and act as if they were processing the
 9202 concatenation of those files. For example:

```
9203        asa file1 file2
```

9204 and:

```
9205        cat file1 file2 | asa
```

9206 have similar results when questions of file access, errors, and performance are ignored. Other
 9207 utilities such as *grep* or *wc* have completely different results in these two cases. This latter type of
 9208 utility is always identified in its DESCRIPTION or OPERANDS sections, whereas the former is
 9209 not. Although it might be possible to create a general assertion about the former case, the
 9210 following points must be addressed:

- 9211 • Access times for the files might be different in the operand case *versus* the *cat* case.
- 9212 • The utility may have error messages that are cognizant of the input filename, and this added
 9213 value should not be suppressed. (As an example, *awk* sets a variable with the filename at
 9214 each file boundary.)

9215 **STDIN**

9216 There is no additional rationale provided for this section.

9217 **INPUT FILES**

9218 A conforming application cannot assume the following three commands are equivalent:

```
9219        tail -n +2 file  

  9220        (sed -n 1q; cat) < file  

  9221        cat file | (sed -n 1q; cat)
```

9222 The second command is equivalent to the first only when the file is seekable. In the third
 9223 command, if the file offset in the open file description were not unspecified, *sed* would have to be
 9224 implemented so that it read from the pipe 1 byte at a time or it would have to employ some
 9225 method to seek backwards on the pipe. Such functionality is not defined currently in POSIX.1
 9226 and does not exist on all historical systems. Other utilities, such as *head*, *read*, and *sh*, have similar
 9227 properties, so the restriction is described globally in this section.

9228 The definition of *text file* is strictly enforced for input to the standard utilities; very few of them
 9229 list exceptions to the undefined results called for here. (Of course, “undefined” here does not
 9230 mean that historical implementations necessarily have to change to start indicating error
 9231 conditions. Conforming applications cannot rely on implementations succeeding or failing when
 9232 non-text files are used.)

9233 The utilities that allow line continuation are generally those that accept input languages, rather
 9234 than pure data. It would be unusual for an input line of this type to exceed `{LINE_MAX}` bytes
 9235 and unreasonable to require that the implementation allow unlimited accumulation of multiple
 9236 lines, each of which could reach `{LINE_MAX}`. Thus, for a conforming application the total of all
 9237 the continued lines in a set cannot exceed `{LINE_MAX}`.

9238 The format description is intended to be sufficiently rigorous to allow other applications to
9239 generate these input files. However, since <blank>s can legitimately be included in some of the
9240 fields described by the standard utilities, particularly in locales other than the POSIX locale, this
9241 intent is not always realized.

9242 **ENVIRONMENT VARIABLES**

9243 There is no additional rationale provided for this section.

9244 **ASYNCHRONOUS EVENTS**

9245 Because there is no language prohibiting it, a utility is permitted to catch a signal, perform some
9246 additional processing (such as deleting temporary files), restore the default signal action (or
9247 action inherited from the parent process), and resignal itself.

9248 **STDOUT**

9249 The format description is intended to be sufficiently rigorous to allow post-processing of output
9250 by other programs, particularly by an *awk* or *lex* parser.

9251 **STDERR**

9252 This section does not describe error messages that refer to incorrect operation of the utility.
9253 Consider a utility that processes program source code as its input. This section is used to
9254 describe messages produced by a correctly operating utility that encounters an error in the
9255 program source code on which it is processing. However, a message indicating that the utility
9256 had insufficient memory in which to operate would not be described.

9257 Some utilities have traditionally produced warning messages without returning a non-zero exit
9258 status; these are specifically noted in their sections. Other utilities shall not write to standard
9259 error if they complete successfully, unless the implementation provides some sort of extension
9260 to increase the verbosity or debugging level.

9261 The format descriptions are intended to be sufficiently rigorous to allow post-processing of
9262 output by other programs.

9263 **OUTPUT FILES**

9264 The format description is intended to be sufficiently rigorous to allow post-processing of output
9265 by other programs, particularly by an *awk* or *lex* parser.

9266 Receipt of the SIGQUIT signal should generally cause termination (unless in some debugging
9267 mode) that would bypass any attempted recovery actions.

9268 **EXTENDED DESCRIPTION**

9269 There is no additional rationale provided for this section.

9270 **EXIT STATUS**

9271 Note the additional discussion of exit values in *Exit Status for Commands* in the *sh* utility. It
9272 describes requirements for returning exit values greater than 125.

9273 A utility may list zero as a successful return, 1 as a failure for a specific reason, and greater than
9274 1 as “an error occurred”. In this case, unspecified conditions may cause a 2 or 3, or other value,
9275 to be returned. A strictly conforming application should be written so that it tests for successful
9276 exit status values (zero in this case), rather than relying upon the single specific error value listed
9277 in IEEE Std 1003.1-200x. In that way, it will have maximum portability, even on implementations

- 9278 with extensions.
- 9279 The standard developers are aware that the general non-enumeration of errors makes it difficult
9280 to write test suites that test the *incorrect* operation of utilities. There are some historical
9281 implementations that have expended effort to provide detailed status messages and a helpful
9282 environment to bypass or explain errors, such as prompting, retrying, or ignoring unimportant
9283 syntax errors; other implementations have not. Since there is no realistic way to mandate system
9284 behavior in cases of undefined application actions or system problems—in a manner acceptable
9285 to all cultures and environments—attention has been limited to the correct operation of utilities
9286 by the conforming application. Furthermore, the conforming application does not need detailed
9287 information concerning errors that it caused through incorrect usage or that it cannot correct.
- 9288 There is no description of defaults for this section because all of the standard utilities specify
9289 something (or explicitly state “Unspecified”) for exit status.
- 9290 **CONSEQUENCES OF ERRORS**
- 9291 Several actions are possible when a utility encounters an error condition, depending on the
9292 severity of the error and the state of the utility. Included in the possible actions of various
9293 utilities are: deletion of temporary or intermediate work files; deletion of incomplete files; and
9294 validity checking of the file system or directory.
- 9295 The text about recursive traversing is meant to ensure that utilities such as *find* process as many
9296 files in the hierarchy as they can. They should not abandon all of the hierarchy at the first error
9297 and resume with the next command-line operand, but should attempt to keep going.
- 9298 **APPLICATION USAGE**
- 9299 This section provides additional caveats, issues, and recommendations to the developer.
- 9300 **EXAMPLES**
- 9301 This section provides sample usage.
- 9302 **RATIONALE**
- 9303 There is no additional rationale provided for this section.
- 9304 **FUTURE DIRECTIONS**
- 9305 FUTURE DIRECTIONS sections act as pointers to related work that may impact the interface in
9306 the future, and often cautions the developer to architect the code to account for a change in this
9307 area. Note that a future directions statement should not be taken as a commitment to adopt a
9308 feature or interface in the future.
- 9309 **SEE ALSO**
- 9310 There is no additional rationale provided for this section.

9311 **CHANGE HISTORY**

9312 There is no additional rationale provided for this section.

9313 **C.1.12 Considerations for Utilities in Support of Files of Arbitrary Size**

9314 This section is intended to clarify the requirements for utilities in support of large files.

9315 The utilities listed in this section are utilities which are used to perform administrative tasks
 9316 such as to create, move, copy, remove, change the permissions, or measure the resources of a
 9317 file. They are useful both as end-user tools and as utilities invoked by applications during
 9318 software installation and operation.

9319 The *chgrp*, *chmod*, *chown*, *ln*, and *rm* utilities probably require use of large file capable versions of
 9320 *stat()*, *lstat()*, *ftw()*, and the **stat** structure.

9321 The *cat*, *cksum*, *cmp*, *cp*, *dd*, *mv*, *sum*, and *touch* utilities probably require use of large file capable
 9322 versions of *creat()*, *open()*, and *fopen()*.

9323 The *cat*, *cksum*, *cmp*, *dd*, *df*, *du*, *ls*, and *sum* utilities may require writing large integer values. For
 9324 example:

- 9325 • The *cat* utility might have a **-n** option which counts <newline>s.
- 9326 • The *cksum* and *ls* utilities report file sizes.
- 9327 • The *cmp* utility reports the line number at which the first difference occurs, and also has a **-l**
 9328 option which reports file offsets.
- 9329 • The *dd*, *df*, *du*, *ls*, and *sum* utilities report block counts.

9330 The *dd*, *find*, and *test* utilities may need to interpret command arguments that contain 64-bit
 9331 values. For *dd*, the arguments include *skip=n*, *seek=n*, and *count=n*. For *find*, the arguments
 9332 include **-sizen**. For *test*, the arguments are those associated with algebraic comparisons.

9333 The *df* utility might need to access large file systems with *statvfs()*.

9334 The *ulimit* utility will need to use large file capable versions of *getrlimit()* and *setrlimit()* and be
 9335 able to read and write large integer values.

9336 **C.1.13 Built-In Utilities**

9337 All of these utilities can be *exec*-ed. There is no requirement that these utilities are actually built
 9338 into the shell itself, but many shells need the capability to do so because the Shell and Utilities
 9339 volume of IEEE Std 1003.1-200x, Section 2.9.1.1, Command Search and Execution requires that
 9340 they be found prior to the *PATH* search. The shell could satisfy its requirements by keeping a list
 9341 of the names and directly accessing the file-system versions regardless of *PATH*. Providing all of
 9342 the required functionality for those such as *cd* or *read* would be more difficult.

9343 There were originally three justifications for allowing the omission of *exec*-able versions:

- 9344 1. It would require wasting space in the file system, at the expense of very small systems.
 9345 However, it has been pointed out that all 16 utilities in the table can be provided with 16
 9346 links to a single-line shell script:
 9347 \$0 "\$@"
- 9348 2. It is not logical to require invocation of utilities such as *cd* because they have no value
 9349 outside the shell environment or cannot be useful in a child process. However, counter-
 9350 examples always seemed to be available for even the most unusual cases:

9351 `find . -type d -exec cd {} \; -exec foo {} \;` |
 9352 (which invokes “foo” on accessible directories) |

9353 `ps ... | sed ... | xargs kill` |

9354 `find . -exec true \; -a ...` |
 9355 (where “true” is used for temporary debugging) |

9356 3. It is confusing to have a utility such as *kill* that can easily be in the file system in the base |
 9357 standard, but that requires built-in status for the UPE (for the % job control job ID notation). |
 9358 It was decided that it was more appropriate to describe the required functionality (rather |
 9359 than the implementation) to the system implementors and let them decide how to satisfy |
 9360 it. |

9361 On the other hand, it was realized that any distinction like this between utilities was not useful |
 9362 to applications, and that the cost to correct it was small. These arguments were ultimately the |
 9363 most effective. |

9364 There were varying reasons for including utilities in the table of built-ins: |

9365 *alias, fc, unalias* |
 9366 The functionality of these utilities is performed more simply within the shell itself and that |
 9367 is the model most historical implementations have used. |

9368 *bg, fg, jobs* |
 9369 All of the job control-related utilities are eligible for built-in status because that is the model |
 9370 most historical implementations have used. |

9371 *cd, getopts, newgrp, read, umask, wait* |
 9372 The functionality of these utilities is performed more simply within the context of the |
 9373 current process. An example can be taken from the usage of the *cd* utility. The purpose of |
 9374 the utility is to change the working directory for subsequent operations. The actions of *cd* |
 9375 affect the process in which *cd* is executed and all subsequent child processes of that process. |
 9376 Based on the ISO POSIX-1 standard p1 process model, changes in the process environment |
 9377 of a child process have no effect on the parent process. If the *cd* utility were executed from a |
 9378 child process, the working directory change would be effective only in the child process. |
 9379 Child processes initiated subsequent to the child process that executed the *cd* utility would |
 9380 not have a changed working directory relative to the parent process. |

9381 *command* |
 9382 This utility was placed in the table primarily to protect scripts that are concerned about |
 9383 their *PATH* being manipulated. The “secure” shell script example in the *command* utility in |
 9384 the Shell and Utilities volume of IEEE Std 1003.1-200x would not be possible if a *PATH* |
 9385 change retrieved an alien version of *command*. (An alternative would have been to |
 9386 implement *getconf* as a built-in, but the standard developers considered that it carried too |
 9387 many changing configuration strings to require in the shell.) |

9388 *kill* |
 9389 Since *kill* provides optional job control functionality using shell notation (%1, %2, and so on), |
 9390 some implementations would find it extremely difficult to provide this outside the shell. |

9391 *true, false* |
 9392 These are in the table as a courtesy to programmers who wish to use the “**while true**” shell |
 9393 construct without protecting *true* from *PATH* searches. (It is acknowledged that “**while\:**” |
 9394 also works, but the idiom with *true* is historically pervasive.) |

9395 All utilities, including those in the table, are accessible via the *system()* and *popen()* functions in |
 9396 the System Interfaces volume of IEEE Std 1003.1-200x. There are situations where the return |

9397 functionality of *system()* and *popen()* is not desirable. Applications that require the exit status of |
 9398 the invoked utility will not be able to use *system()* or *popen()*, since the exit status returned is |
 9399 that of the command language interpreter rather than that of the invoked utility. The alternative |
 9400 for such applications is the use of the *exec* family. |

9401 **C.2 Shell Command Language**

9402 **C.2.1 Shell Introduction**

9403 The System V shell was selected as the starting point for the Shell and Utilities volume of
 9404 IEEE Std 1003.1-200x. The BSD C shell was excluded from consideration for the following
 9405 reasons:

- 9406 • Most historically portable shell scripts assume the Version 7 Bourne shell, from which the
 9407 System V shell is derived.
- 9408 • The majority of tutorial materials on shell programming assume the System V shell.

9409 The construct "#!" is reserved for implementations wishing to provide that extension. If it were
 9410 not reserved, the Shell and Utilities volume of IEEE Std 1003.1-200x would disallow it by forcing
 9411 it to be a comment. As it stands, a strictly conforming application must not use "#!" as the first
 9412 two characters of the file.

9413 **C.2.2 Quoting**

9414 There is no additional rationale provided for this section.

9415 *C.2.2.1 Escape Character (Backslash)*

9416 There is no additional rationale provided for this section.

9417 *C.2.2.2 Single-Quotes*

9418 A backslash cannot be used to escape a single-quote in a single-quoted string. An embedded
 9419 quote can be created by writing, for example: "'a'\ 'b'", which yields "a'b". (See the Shell
 9420 and Utilities volume of IEEE Std 1003.1-200x, Section 2.6.5, Field Splitting for a better
 9421 understanding of how portions of words are either split into fields or remain concatenated.) A
 9422 single token can be made up of concatenated partial strings containing all three kinds of quoting
 9423 or escaping, thus permitting any combination of characters.

9424 *C.2.2.3 Double-Quotes*

9425 The escaped <newline> used for line continuation is removed entirely from the input and is not
 9426 replaced by any white space. Therefore, it cannot serve as a token separator.

9427 In double-quoting, if a backslash is immediately followed by a character that would be
 9428 interpreted as having a special meaning, the backslash is deleted and the subsequent character is
 9429 taken literally. If a backslash does not precede a character that would have a special meaning, it
 9430 is left in place unmodified and the character immediately following it is also left unmodified.
 9431 Thus, for example:

9432 "\\$" → \$ |

9433 "\a" → \a |

9434 It would be desirable to include the statement “The characters from an enclosed “\${” to the
 9435 matching ’}’ shall not be affected by the double quotes”, similar to the one for “\$()”.
 9436 However, historical practice in the System V shell prevents this.

9437 The requirement that double-quotes be matched inside “\${ . . . }” within double-quotes and the
 9438 rule for finding the matching ’}’ in the Shell and Utilities volume of IEEE Std 1003.1-200x,
 9439 Section 2.6.2, Parameter Expansion eliminate several subtle inconsistencies in expansion for
 9440 historical shells in rare cases; for example:

```
9441     "${foo-bar}"
```

9442 yields **bar** when **foo** is not defined, and is an invalid substitution when **foo** is defined, in many
 9443 historical shells. The differences in processing the “\${ . . . }” form have led to inconsistencies
 9444 between historical systems. A consequence of this rule is that single-quotes cannot be used to
 9445 quote the ’}’ within “\${ . . . }”; for example:

```
9446     unset bar
9447     foo="${bar-'}'"
```

9448 is invalid because the “\${ . . . }” substitution contains an unpaired unescaped single-quote. The
 9449 backslash can be used to escape the ’}’ in this example to achieve the desired result:

```
9450     unset bar
9451     foo="${bar-\}]"
```

9452 The differences in processing the “\${ . . . }” form have led to inconsistencies between the
 9453 historical System V shell, BSD, and KornShells, and the text in the Shell and Utilities volume of
 9454 IEEE Std 1003.1-200x is an attempt to converge them without breaking too many applications.
 9455 The only alternative to this compromise between shells would be to make the behavior
 9456 unspecified whenever the literal characters ‘ ’, ‘{’, ‘}’, and ‘”’ appear within “\${ . . . }”.
 9457 To write a portable script that uses these values, a user would have to assign variables; for
 9458 example:

```
9459     squote=\' dquote=\" lbrace='{ rbrace=}'
9460     ${foo-$$squote$$rbrace$$squote}
```

9461 rather than:

```
9462     ${foo-"' }'"}
```

9463 Some implementations have allowed the end of the word to terminate the backquoted command
 9464 substitution, such as in:

```
9465     "`echo hello"
```

9466 This usage is undefined; the matching backquote is required by the Shell and Utilities volume of
 9467 IEEE Std 1003.1-200x. The other undefined usage can be illustrated by the example:

```
9468     sh -c '` echo "foo`'
```

9469 The description of the recursive actions involving command substitution can be illustrated with
 9470 an example. Upon recognizing the introduction of command substitution, the shell parses input
 9471 (in a new context), gathering the source for the command substitution until an unbalanced ’)’
 9472 or ‘`’ is located. For example, in the following:

```
9473     echo "$(date; echo "  

  9474         one" )"
```

9475 the double-quote following the *echo* does not terminate the first double-quote; it is part of the
 9476 command substitution script. Similarly, in:

9477 echo "\$ (echo *)" |
 9478 the asterisk is not quoted since it is inside command substitution; however:
 9479 echo "\$ (echo "*")" |
 9480 is quoted (and represents the asterisk character itself).

9481 C.2.3 Token Recognition

9482 The "(" and ")" symbols are control operators in the KornShell, used for an alternative |
 9483 syntax of an arithmetic expression command. A conforming application cannot use "(" as a |
 9484 single token (with the exception of the "\$ (" form for shell arithmetic).

9485 On some implementations, the symbol "(" is a control operator; its use produces unspecified |
 9486 results. Applications that wish to have nested subshells, such as:

```
9487           ((echo Hello);(echo World)) |
```

9488 shall separate the "(" characters into two tokens by including white space between them. |
 9489 Some systems may treat these as invalid arithmetic expressions instead of subshells.

9490 Certain combinations of characters are invalid in portable scripts, as shown in the grammar. |
 9491 Implementations may use these combinations (such as "|&") as valid control operators. Portable |
 9492 scripts cannot rely on receiving errors in all cases where this volume of IEEE Std 1003.1-200x |
 9493 indicates that a syntax is invalid.

9494 The (3) rule about combining characters to form operators is not meant to preclude systems from |
 9495 extending the shell language when characters are combined in otherwise invalid ways. |
 9496 Conforming applications cannot use invalid combinations, and test suites should not penalize |
 9497 systems that take advantage of this fact. For example, the unquoted combination "|&" is not |
 9498 valid in a POSIX script, but has a specific KornShell meaning.

9499 The (10) rule about '#' as the current character is the first in the sequence in which a new token |
 9500 is being assembled. The '#' starts a comment only when it is at the beginning of a token. This |
 9501 rule is also written to indicate that the search for the end-of-comment does not consider escaped |
 9502 <newline> specially, so that a comment cannot be continued to the next line.

9503 C.2.3.1 Alias Substitution

9504 The alias capability was added in the UPE because it is widely used in historical |
 9505 implementations by interactive users.

9506 The definition of *alias name* precludes an alias name containing a slash character. Since the text |
 9507 applies to the command words of simple commands, reserved words (in their proper places) |
 9508 cannot be confused with aliases.

9509 The placement of alias substitution in token recognition makes it clear that it precedes all of the |
 9510 word expansion steps.

9511 An example concerning trailing <blank>s and reserved words follows. If the user types:

```
9512           $ alias foo="/bin/ls " |  

  9513           $ alias while="/ " |
```

9514 The effect of executing:

```

9515     $ while true
9516     > do
9517     > echo "Hello, World"
9518     > done

```

9519 is a never-ending sequence of "Hello, World" strings to the screen. However, if the user
 9520 types:

```

9521     $ foo while

```

9522 the result is an *ls* listing of /. Since the alias substitution for **foo** ends in a <space>, the next word
 9523 is checked for alias substitution. The next word, **while**, has also been aliased, so it is substituted
 9524 as well. Since it is not in the proper position as a command word, it is not recognized as a
 9525 reserved word.

9526 If the user types:

```

9527     $ foo; while

```

9528 **while** retains its normal reserved-word properties.

9529 C.2.4 Reserved Words

9530 All reserved words are recognized syntactically as such in the contexts described. However, note
 9531 that **in** is the only meaningful reserved word after a **case** or **for**; similarly, **in** is not meaningful as
 9532 the first word of a simple command.

9533 Reserved words are recognized only when they are delimited (that is, meet the definition of the
 9534 Base Definitions volume of IEEE Std 1003.1-200x, Section 3.435, Word), whereas operators are
 9535 themselves delimiters. For instance, '(' and ') ' are control operators, so that no <space> is
 9536 needed in (*list*). However, '{ ' and ' } ' are reserved words in { *list*; }, so that in this case the
 9537 leading <space> and semicolon are required.

9538 The list of unspecified reserved words is from the KornShell, so conforming applications cannot
 9539 use them in places a reserved word would be recognized. This list contained **time** in early
 9540 proposals, but it was removed when the *time* utility was selected for the Shell and Utilities
 9541 volume of IEEE Std 1003.1-200x.

9542 There was a strong argument for promoting braces to operators (instead of reserved words), so
 9543 they would be syntactically equivalent to subshell operators. Concerns about compatibility
 9544 outweighed the advantages of this approach. Nevertheless, conforming applications should
 9545 consider quoting '{ ' and ' } ' when they represent themselves.

9546 The restriction on ending a name with a colon is to allow future implementations that support
 9547 named labels for flow control; see the RATIONALE for the *break* built-in utility .

9548 It is possible that a future version of the Shell and Utilities volume of IEEE Std 1003.1-200x may
 9549 require that '{ ' and ' } ' be treated individually as control operators, although the token "{ }"
 9550 will probably be a special-case exemption from this because of the often-used *find*{ } construct.

9551 **C.2.5 Parameters and Variables**9552 *C.2.5.1 Positional Parameters*

9553 There is no additional rationale provided for this section.

9554 *C.2.5.2 Special Parameters*

9555 Most historical implementations implement subshells by forking; thus, the special parameter
 9556 ' \$ ' does not necessarily represent the process ID of the shell process executing the commands
 9557 since the subshell execution environment preserves the value of ' \$ '.

9558 If a subshell were to execute a background command, the value of "\$!" for the parent would
 9559 not change. For example:

```
9560 (
9561   date &
9562   echo $!
9563 )
9564 echo $!
```

9565 would echo two different values for "\$!".

9566 The "\$-" special parameter can be used to save and restore *set* options:

```
9567   Save=$(echo $- | sed 's/[ics]//g')
9568   ...
9569   set +aCefnuvx
9570   if [ -n "$Save" ]; then
9571       set -$Save
9572   fi
```

9573 The three options are removed using *sed* in the example because they may appear in the value of
 9574 "\$-" (from the *sh* command line), but are not valid options to *set*.

9575 The descriptions of parameters '*' and '@' assume the reader is familiar with the field splitting
 9576 discussion in the Shell and Utilities volume of IEEE Std 1003.1-200x, Section 2.6.5, Field Splitting
 9577 and understands that portions of the word remain concatenated unless there is some reason to
 9578 split them into separate fields.

9579 Some examples of the '*' and '@' properties, including the concatenation aspects:

```
9580   set "abc" "def ghi" "jkl"
9581   echo $*      => "abc" "def" "ghi" "jkl"
9582   echo "$*"   => "abc def ghi jkl"
9583   echo $@     => "abc" "def" "ghi" "jkl"
```

9584 **but:**

```
9585   echo "$@"    => "abc" "def ghi" "jkl"
9586   echo "xx$@yy" => "xxabc" "def ghi" "jkl"yy"
9587   echo "$@$@"  => "abc" "def ghi" "jklabc" "def ghi" "jkl"
```

9588 In the preceding examples, the double-quote characters that appear after the "=>" do not appear
 9589 in the output and are used only to illustrate word boundaries.

9590 The following example illustrates the effect of setting *IFS* to a null string:

```

9591      $ IFS=' '
9592      $ set foo bar bam
9593      $ echo "$@"
9594      foo bar bam
9595      $ echo "$*"
9596      foobarbam
9597      $ unset IFS
9598      $ echo "$*"
9599      foo bar bam

```

9600 C.2.5.3 Shell Variables

9601 See the discussion of *IFS* in Section C.2.6.5 (on page 3529).

9602 The prohibition on *LC_CTYPE* changes affecting lexical processing protects the shell
 9603 implementor (and the shell programmer) from the ill effects of changing the definition of
 9604 <blank> or the set of alphabetic characters in the current environment. It would probably not be
 9605 feasible to write a compiled version of a shell script without this rule. The rule applies only to
 9606 the current invocation of the shell and its subshells—invoking a shell script or performing *exec sh*
 9607 would subject the new shell to the changes in *LC_CTYPE*.

9608 Other common environment variables used by historical shells are not specified by the Shell and
 9609 Utilities volume of IEEE Std 1003.1-200x, but they should be reserved for the historical uses.

9610 Tilde expansion for components of the *PATH* in an assignment such as:

```
9611     PATH=~hlj/bin:~dwc/bin:$PATH
```

9612 is a feature of some historical shells and is allowed by the wording of the Shell and Utilities
 9613 volume of IEEE Std 1003.1-200x, Section 2.6.1, Tilde Expansion. Note that the tildes are expanded
 9614 during the assignment to *PATH*, not when *PATH* is accessed during command search.

9615 The following entries represent additional information about variables included in the Shell and
 9616 Utilities volume of IEEE Std 1003.1-200x, or rationale for common variables in use by shells that
 9617 have been excluded:

9618 — (Underscore.) While underscore is historical practice, its overloaded usage in
 9619 the KornShell is confusing, and it has been omitted from the Shell and Utilities
 9620 volume of IEEE Std 1003.1-200x.

9621 *ENV* This variable can be used to set aliases and other items local to the invocation
 9622 of a shell. The file referred to by *ENV* differs from *\$HOME/.profile* in that
 9623 *.profile* is typically executed at session start-up, whereas the *ENV* file is
 9624 executed at the beginning of each shell invocation. The *ENV* value is
 9625 interpreted in a manner similar to a dot script, in that the commands are
 9626 executed in the current environment and the file needs to be readable, but not
 9627 executable. However, unlike dot scripts, no *PATH* searching is performed.
 9628 This is used as a guard against Trojan Horse security breaches.

9629 *ERRNO* This variable was omitted from the Shell and Utilities volume of
 9630 IEEE Std 1003.1-200x because the values of error numbers are not defined in
 9631 IEEE Std 1003.1-200x in a portable manner.

9632 *FCEDIT* Since this variable affects only the *fc* utility, it has been omitted from this more
 9633 global place. The value of *FCEDIT* does not affect the command line editing
 9634 mode in the shell; see the description of *set -o vi* in the *set* built-in utility.

9635	<i>PS1</i>	This variable is used for interactive prompts. Historically, the “superuser”
9636		has had a prompt of '#'. Since privileges are not required to be monolithic, it
9637		is difficult to define which privileges should cause the alternate prompt.
9638		However, a sufficiently powerful user should be reminded of that power by
9639		having an alternate prompt.
9640	<i>PS3</i>	This variable is used by the KornShell for the <i>select</i> command. Since the POSIX
9641		shell does not include <i>select</i> , <i>PS3</i> was omitted.
9642	<i>PS4</i>	This variable is used for shell debugging. For example, the following script:
9643		<pre>PS4='[\${LINENO}]+ '</pre>
9644		<pre>set -x</pre>
9645		<pre>echo Hello</pre>
9646		writes the following to standard error:
9647		<pre>[3]+ echo Hello</pre>
9648	<i>RANDOM</i>	This pseudo-random number generator was not seen as being useful to
9649		interactive users.
9650	<i>SECONDS</i>	Although this variable is sometimes used with <i>PS1</i> to allow the display of the
9651		current time in the prompt of the user, it is not one that would be manipulated
9652		frequently enough by an interactive user to include in the Shell and Utilities
9653		volume of IEEE Std 1003.1-200x.

9654 C.2.6 Word Expansions

9655 Step (2) refers to the “portions of fields generated by step (1)”. For example, if the word being
 9656 expanded were "\$x+\$y" and *IFS*=+, the word would be split only if "\$x" or "\$y" contained
 9657 '+'; the '+' in the original word was not generated by step (1).

9658 *IFS* is used for performing field splitting on the results of parameter and command substitution;
 9659 it is not used for splitting all fields. Previous versions of the shell used it for splitting all fields
 9660 during field splitting, but this has severe problems because the shell can no longer parse its own
 9661 script. There are also important security implications caused by this behavior. All useful
 9662 applications of *IFS* use it for parsing input of the *read* utility and for splitting the results of
 9663 parameter and command substitution.

9664 The rule concerning expansion to a single field requires that if **foo=abc** and **bar=def**, that:

```
9665 "$foo" "$bar"
```

9666 expands to the single field:

```
9667 abcdef
```

9668 The rule concerning empty fields can be illustrated by:

```

9669      $   unset foo
9670      $   set $foo bar ' ' xyz "$foo" abc
9671      $   for i
9672      >   do
9673      >     echo "-$i-"
9674      >   done
9675      -bar-
9676      --
9677      -xyz-
9678      --
9679      -abc-

```

9680 Step (1) indicates that parameter expansion, command substitution, and arithmetic expansion
 9681 are all processed simultaneously as they are scanned. For example, the following is valid
 9682 arithmetic:

```

9683      x=1
9684      echo $(( $(echo 3)+$x ))

```

9685 An early proposal stated that tilde expansion preceded the other steps, but this is not the case in
 9686 known historical implementations; if it were, and if a referenced home directory contained a '\$'
 9687 character, expansions would result within the directory name.

9688 C.2.6.1 Tilde Expansion

9689 Tilde expansion generally occurs only at the beginning of words, but an exception based on
 9690 historical practice has been included:

```

9691      PATH=/posix/bin:~djk/bin

```

9692 This is eligible for tilde expansion because tilde follows a colon and none of the relevant
 9693 characters is quoted. Consideration was given to prohibiting this behavior because any of the
 9694 following are reasonable substitutes:

```

9695      PATH=$(printf %s ~karels/bin : ~bostic/bin)
9696      for Dir in ~maat/bin ~srb/bin ...
9697      do
9698          PATH=${PATH:+$PATH:}$Dir
9699      done

```

9700 In the first command, explicit colons are used for each directory. In all cases, the shell performs
 9701 tilde expansion on each directory because all are separate words to the shell.

9702 Note that expressions in operands such as:

```

9703      make -k mumble LIBDIR=~chet/lib

```

9704 do not qualify as shell variable assignments, and tilde expansion is not performed (unless the
 9705 command does so itself, which *make* does not).

9706 Because of the requirement that the word is not quoted, the following are not equivalent; only
 9707 the last causes tilde expansion:

```

9708      \~hlj/   ~h\lj/   ~"hlj"/   ~hlj\   ~hlj/

```

9709 In an early proposal, tilde expansion occurred following any unquoted equals sign or colon, but
 9710 this was removed because of its complexity and to avoid breaking commands such as:

9711 rcp hostname:~marc/.profile .

9712 A suggestion was made that the special sequence "\$~" should be allowed to force tilde
 9713 expansion anywhere. Since this is not historical practice, it has been left for future
 9714 implementations to evaluate. (The description in the Shell and Utilities volume of
 9715 IEEE Std 1003.1-200x, Section 2.2, Quoting requires that a dollar sign be quoted to represent
 9716 itself, so the "\$~" combination is already unspecified.)

9717 The results of giving tilde with an unknown login name are undefined because the KornShell
 9718 "~+" and "~-" constructs make use of this condition, but in general it is an error to give an
 9719 incorrect login name with tilde. The results of having *HOME* unset are unspecified because some
 9720 historical shells treat this as an error.

9721 C.2.6.2 Parameter Expansion

9722 The rule for finding the closing '}' in "\${...}" is the one used in the KornShell and is
 9723 upwardly-compatible with the Bourne shell, which does not determine the closing '}' until the
 9724 word is expanded. The advantage of this is that incomplete expansions, such as:

9725 \${foo

9726 can be determined during tokenization, rather than during expansion.

9727 The string length and substring capabilities were included because of the demonstrated need for
 9728 them, based on their usage in other shells, such as C shell and KornShell.

9729 Historical versions of the KornShell have not performed tilde expansion on the word part of
 9730 parameter expansion; however, it is more consistent to do so.

9731 C.2.6.3 Command Substitution

9732 The "\$()" form of command substitution solves a problem of inconsistent behavior when using
 9733 backquotes. For example:

Command	Output
echo '\\$x'	\\$x
echo `echo '\\$x'`	\$x
echo \$(echo '\\$x')	\\$x

9738 Additionally, the backquoted syntax has historical restrictions on the contents of the embedded
 9739 command. While the newer "\$()" form can process any kind of valid embedded script, the
 9740 backquoted form cannot handle some valid scripts that include backquotes. For example, these
 9741 otherwise valid embedded scripts do not work in the left column, but do work on the right:

9742 echo `	echo \$(
9743 cat <<\eof	cat <<\eof
9744 a here-doc with `	a here-doc with)
9745 eof	eof
9746 `)
9747 echo `	echo \$(
9748 echo abc # a comment with `	echo abc # a comment with)
9749 `)
9750 echo `	echo \$(
9751 echo ` ` `	echo ` ` `
9752 `)

9753 Because of these inconsistent behaviors, the backquoted variety of command substitution is not
 9754 recommended for new applications that nest command substitutions or attempt to embed
 9755 complex scripts.

9756 The KornShell feature:

9757 If *command* is of the form `<word`, *word* is expanded to generate a pathname, and the value of |
 9758 the command substitution is the contents of this file with any trailing `<newline>`s deleted. |

9759 was omitted from the Shell and Utilities volume of IEEE Std 1003.1-200x because `$(cat word)` is
 9760 an appropriate substitute. However, to prevent breaking numerous scripts relying on this
 9761 feature, it is unspecified to have a script within `"$()"` that has only redirections.

9762 The requirement to separate `"$("` and `'('` when a single subshell is command-substituted is to
 9763 avoid any ambiguities with arithmetic expansion.

9764 C.2.6.4 Arithmetic Expansion

9765 The `"()"` form of KornShell arithmetic in early proposals was omitted. The standard
 9766 developers concluded that there was a strong desire for some kind of arithmetic evaluator to
 9767 replace *expr*, and that relating it to `'$'` makes it work well with the standard shell language, and
 9768 it provides access to arithmetic evaluation in places where accessing a utility would be
 9769 inconvenient.

9770 The syntax and semantics for arithmetic were changed for the ISO/IEC 9945-2:1993 standard.
 9771 The language is essentially a pure arithmetic evaluator of constants and operators (excluding
 9772 assignment) and represents a simple subset of the previous arithmetic language (which was
 9773 derived from the KornShell `"()"` construct). The syntax was changed from that of a
 9774 command denoted by `((expression))` to an expansion denoted by `$(expression)`. The new form is
 9775 a dollar expansion `'$'` that evaluates the expression and substitutes the resulting value.
 9776 Objections to the previous style of arithmetic included that it was too complicated, did not fit in
 9777 well with the use of variables in the shell, and its syntax conflicted with subshells. The
 9778 justification for the new syntax is that the shell is traditionally a macro language, and if a new
 9779 feature is to be added, it should be accomplished by extending the capabilities presented by the
 9780 current model of the shell, rather than by inventing a new one outside the model; adding a new
 9781 dollar expansion was perceived to be the most intuitive and least destructive way to add such a
 9782 new capability.

9783 In early proposals, a form `$(expression)` was used. It was functionally equivalent to the `"$()"`
 9784 of the current text, but objections were lodged that the 1988 KornShell had already implemented
 9785 `"$()"` and there was no compelling reason to invent yet another syntax. Furthermore, the
 9786 `"$[]"` syntax had a minor incompatibility involving the patterns in `case` statements.

9787 The portion of the ISO C standard arithmetic operations selected corresponds to the operations
 9788 historically supported in the KornShell.

9789 It was concluded that the `test` command (`D`) was sufficient for the majority of relational arithmetic
 9790 tests, and that tests involving complicated relational expressions within the shell are rare, yet
 9791 could still be accommodated by testing the value of `"$()"` itself. For example:

```
9792     # a complicated relational expression
9793     while [ $( (( $x + $y ) / ( $a * $b )) < ( $foo * $bar )) -ne 0 ]
```

9794 or better yet, the rare script that has many complex relational expressions could define a
 9795 function like this:

```

9796     val() {
9797         return $(!$1)
9798     }

```

9799 and complicated tests would be less intimidating:

```

9800     while val $(( (($x + $y)/($a * $b)) < ($foo*$bar) ))
9801     do
9802         # some calculations
9803     done

```

9804 A suggestion that was not adopted was to modify *true* and *false* to take an optional argument,
 9805 and *true* would exit true only if the argument was non-zero, and *false* would exit false only if the
 9806 argument was non-zero:

```

9807     while true $(( $x > 5 && $y <= 25 ))

```

9808 There is a minor portability concern with the new syntax. The example `$(2+2)` could have been
 9809 intended to mean a command substitution of a utility named `2+2` in a subshell. The standard
 9810 developers considered this to be obscure and isolated to some KornShell scripts (because `"$()"`
 9811 command substitution existed previously only in the KornShell). The text on command
 9812 substitution requires that the `"$("` and `'('` be separate tokens if this usage is needed.

9813 An example such as:

```

9814     echo $((echo hi);(echo there))

```

9815 should not be misinterpreted by the shell as arithmetic because attempts to balance the
 9816 parentheses pairs would indicate that they are subshells. However, as indicated by the Base
 9817 Definitions volume of IEEE Std 1003.1-200x, Section 3.112, Control Operator, a conforming
 9818 application must separate two adjacent parentheses with white space to indicate nested
 9819 subshells.

9820 Although the ISO/IEC 9899:1999 standard now requires support for **long long** and allows
 9821 extended integer types with higher ranks, IEEE Std 1003.1-200x only requires arithmetic
 9822 expansions to support **signed long** integer arithmetic. Implementations are encouraged to
 9823 support signed integer values at least as large as the size of the largest file allowed on the
 9824 implementation.

9825 Implementations are also allowed to perform floating-point evaluations as long as an
 9826 application won't see different results for expressions that would not overflow **signed long**
 9827 integer expression evaluation. (This includes appropriate truncation of results to integer values.)

9828 Changes made in response to IEEE PASC Interpretation 1003.2 #208 removed the requirement
 9829 that the integer constant suffixes `l` and `L` had to be recognized. The ISO POSIX-2: 1993 standard
 9830 didn't require the `u`, `uL`, `uL`, `U`, `UL`, `UL`, `Lu`, `LU`, `Lu`, and `LU` suffixes since only signed integer
 9831 arithmetic was required. Since all arithmetic expressions were treated as handling **signed long**
 9832 integer types anyway, the `l` and `L` suffixes were redundant. No known scripts used them and
 9833 some historic shells didn't support them. When the ISO/IEC 9899: 1999 standard was used as the
 9834 basis for the description of arithmetic processing, the `ll` and `LL` suffixes and combinations were
 9835 also not required. Implementations are still free to accept any or all of these suffices, but are not
 9836 required to do so.

9837 There was also some confusion as to whether the shell was required to recognize character
 9838 constants. Syntactically, character constants were required to be recognized, but the
 9839 requirements for the handling of backslash ("`\\`") and quote ("`'\''`") characters (needed to
 9840 specify character constants) within an arithmetic expansion were ambiguous. Furthermore, no
 9841 known shells supported them. Changes made in response to IEEE PASC Interpretation 1003.2

9842 #208 removed the requirement to support them (if they were indeed required before). |
 9843 IEEE Std 1003.1-200x clearly does not require support for character constants. |

9844 C.2.6.5 Field Splitting

9845 The operation of field splitting using *IFS*, as described in early proposals, was based on the way
 9846 the KornShell splits words, but it is incompatible with other common versions of the shell.
 9847 However, each has merit, and so a decision was made to allow both. If the *IFS* variable is unset
 9848 or is `<space><tab><newline>`, the operation is equivalent to the way the System V shell splits
 9849 words. Using characters outside the `<space><tab><newline>` set yields the KornShell behavior,
 9850 where each of the non-`<space><tab><newline>`s is significant. This behavior, which affords the
 9851 most flexibility, was taken from the way the original *awk* handled field splitting.

9852 Rule (3) can be summarized as a pseudo-ERE:

```
9853 (s*ns*|s+)
```

9854 where *s* is an *IFS* white space character and *n* is a character in the *IFS* that is not white space.
 9855 Any string matching that ERE delimits a field, except that the *s+* form does not delimit fields at
 9856 the beginning or the end of a line. For example, if *IFS* is `<space>/<comma>/<tab>`, the string:

```
9857 <space><space>red<space><space>, <space>white<space>blue
```

9858 yields the three colors as the delimited fields.

9859 C.2.6.6 Pathname Expansion

9860 There is no additional rationale provided for this section.

9861 C.2.6.7 Quote Removal

9862 There is no additional rationale provided for this section.

9863 C.2.7 Redirection

9864 In the System Interfaces volume of IEEE Std 1003.1-200x, file descriptors are integers in the range
 9865 0–(`{OPEN_MAX}–1`). The file descriptors discussed in the Shell and Utilities volume of
 9866 IEEE Std 1003.1-200x, Section 2.7, Redirection are that same set of small integers.

9867 Having multi-digit file descriptor numbers for I/O redirection can cause some obscure
 9868 compatibility problems. Specifically, scripts that depend on an example command:

```
9869 echo 22>/dev/null
```

9870 echoing 2 to standard error or 22 to standard output are no longer portable. However, the file
 9871 descriptor number still must be delimited from the preceding text. For example:

```
9872 cat file2>foo
```

9873 writes the contents of **file2**, not the contents of **file**.

9874 The "`>`" format of output redirection was adopted from the KornShell. Along with the
 9875 *noclobber* option, *set –C*, it provides a safety feature to prevent inadvertent overwriting of
 9876 existing files. (See the RATIONALE for the *pathchk* utility for why this step was taken.) The
 9877 restriction on regular files is historical practice.

9878 The System V shell and the KornShell have differed historically on pathname expansion of *word*;
 9879 the former never performed it, the latter only when the result was a single field (file). As a
 9880 compromise, it was decided that the KornShell functionality was useful, but only as a shorthand
 9881 device for interactive users. No reasonable shell script would be written with a command such

9882 as:

```
9883     cat foo > a*
```

9884 Thus, shell scripts are prohibited from doing it, while interactive users can select the shell with
9885 which they are most comfortable.

9886 The construct `2>&1` is often used to redirect standard error to the same file as standard output.
9887 Since the redirections take place beginning to end, the order of redirections is significant. For
9888 example:

```
9889     ls > foo 2>&1
```

9890 directs both standard output and standard error to file **foo**. However:

```
9891     ls 2>&1 > foo
```

9892 only directs standard output to file **foo** because standard error was duplicated as standard
9893 output before standard output was directed to file **foo**.

9894 The "<>" operator could be useful in writing an application that worked with several terminals,
9895 and occasionally wanted to start up a shell. That shell would in turn be unable to run
9896 applications that run from an ordinary controlling terminal unless it could make use of "<>"
9897 redirection. The specific example is a historical version of the pager *more*, which reads from
9898 standard error to get its commands, so standard input and standard output are both available
9899 for their usual usage. There is no way of saying the following in the shell without "<>":

```
9900     cat food | more - >/dev/tty03 2<>/dev/tty03
```

9901 Another example of "<>" is one that opens `/dev/tty` on file descriptor 3 for reading and writing:

```
9902     exec 3<> /dev/tty
```

9903 An example of creating a lock file for a critical code region:

```
9904     set -C
9905     until    2> /dev/null > lockfile
9906     do      sleep 30
9907     done
9908     set +C
9909     perform critical function
9910     rm lockfile
```

9911 Since `/dev/null` is not a regular file, no error is generated by redirecting to it in *noclobber* mode.

9912 Tilde expansion is not performed on a here-document because the data is treated as if it were
9913 enclosed in double quotes.

9914 C.2.7.1 Redirecting Input

9915 There is no additional rationale provided for this section.

9916 C.2.7.2 Redirecting Output

9917 There is no additional rationale provided for this section.

9918 *C.2.7.3 Appending Redirected Output*

9919 Note that when a file is opened (even with the `O_APPEND` flag set), the initial file offset for that
9920 file is set to the beginning of the file. Some historic shells set the file offset to the current end-of-
9921 file when append mode shell redirection was used, but this is not allowed by
9922 IEEE Std 1003.1-200x.

9923 *C.2.7.4 Here-Document*

9924 There is no additional rationale provided for this section.

9925 *C.2.7.5 Duplicating an Input File Descriptor*

9926 There is no additional rationale provided for this section.

9927 *C.2.7.6 Duplicating an Output File Descriptor*

9928 There is no additional rationale provided for this section.

9929 *C.2.7.7 Open File Descriptors for Reading and Writing*

9930 There is no additional rationale provided for this section.

9931 **C.2.8 Exit Status and Errors**9932 *C.2.8.1 Consequences of Shell Errors*

9933 There is no additional rationale provided for this section.

9934 *C.2.8.2 Exit Status for Commands*

9935 There is a historical difference in *sh* and *ksh* non-interactive error behavior. When a command
9936 named in a script is not found, some implementations of *sh* exit immediately, but *ksh* continues
9937 with the next command. Thus, the Shell and Utilities volume of IEEE Std 1003.1-200x says that
9938 the shell “may” exit in this case. This puts a small burden on the programmer, who has to test
9939 for successful completion following a command if it is important that the next command not be
9940 executed if the previous command was not found. If it is important for the command to have
9941 been found, it was probably also important for it to complete successfully. The test for successful
9942 completion would not need to change.

9943 Historically, shells have returned an exit status of $128+n$, where n represents the signal number.
9944 Since signal numbers are not standardized, there is no portable way to determine which signal
9945 caused the termination. Also, it is possible for a command to exit with a status in the same range
9946 of numbers that the shell would use to report that the command was terminated by a signal.
9947 Implementations are encouraged to choose exit values greater than 256 to indicate programs
9948 that terminate by a signal so that the exit status cannot be confused with an exit status generated
9949 by a normal termination.

9950 Historical shells make the distinction between “utility not found” and “utility found but cannot
9951 execute” in their error messages. By specifying two seldomly used exit status values for these
9952 cases, 127 and 126 respectively, this gives an application the opportunity to make use of this
9953 distinction without having to parse an error message that would probably change from locale to
9954 locale. The *command*, *env*, *nohup*, and *xargs* utilities in the Shell and Utilities volume of
9955 IEEE Std 1003.1-200x have also been specified to use this convention.

9956 When a command fails during word expansion or redirection, most historical implementations
9957 exit with a status of 1. However, there was some sentiment that this value should probably be

9958 much higher so that an application could distinguish this case from the more normal exit status
 9959 values. Thus, the language “greater than zero” was selected to allow either method to be
 9960 implemented.

9961 C.2.9 Shell Commands

9962 A description of an “empty command” was removed from an early proposal because it is only
 9963 relevant in the cases of *sh -c " "*, *system(" ")*, or an empty shell-script file (such as the
 9964 implementation of *true* on some historical systems). Since it is no longer mentioned in the Shell
 9965 and Utilities volume of IEEE Std 1003.1-200x, it falls into the silently unspecified category of
 9966 behavior where implementations can continue to operate as they have historically, but
 9967 conforming applications do not construct empty commands. (However, note that *sh* does
 9968 explicitly state an exit status for an empty string or file.) In an interactive session or a script with
 9969 other commands, extra <newline>s or semicolons, such as;

```
9970     $ false
9971     $
9972     $ echo $?
9973     1
```

9974 would not qualify as the empty command described here because they would be consumed by
 9975 other parts of the grammar.

9976 C.2.9.1 Simple Commands

9977 The enumerated list is used only when the command is actually going to be executed. For
 9978 example, in:

```
9979     true || $foo *
```

9980 no expansions are performed.

9981 The following example illustrates both how a variable assignment without a command name
 9982 affects the current execution environment, and how an assignment with a command name only
 9983 affects the execution environment of the command:

```
9984     $ x=red
9985     $ echo $x
9986     red
9987     $ export x
9988     $ sh -c 'echo $x'
9989     red
9990     $ x=blue sh -c 'echo $x'
9991     blue
9992     $ echo $x
9993     red
```

9994 This next example illustrates that redirections without a command name are still performed:

```
9995     $ ls foo
9996     ls: foo: no such file or directory
9997     $ > foo
9998     $ ls foo
9999     foo
```

10000 A command without a command name, but one that includes a command substitution, has an
 10001 exit status of the last command substitution that the shell performed. For example:

```

10002     if      x=$(command)
10003     then    ...
10004     fi

```

10005 An example of redirections without a command name being performed in a subshell shows that
 10006 the here-document does not disrupt the standard input of the **while** loop:

```

10007     IFS=:
10008     while  read a b
10009     do    echo $a
10010         <<-eof
10011         Hello
10012         eof
10013     done </etc/passwd

```

10014 Some examples of commands without command names in AND-OR lists:

```

10015     > foo || {
10016         echo "error: foo cannot be created" >&2
10017         exit 1
10018     }
10019     # set saved if /vmunix.save exists
10020     test -f /vmunix.save && saved=1

```

10021 Command substitution and redirections without command names both occur in subshells, but
 10022 they are not necessarily the same ones. For example, in:

```

10023     exec 3> file
10024     var=$(echo foo >&3) 3>&1

```

10025 it is unspecified whether **foo** is echoed to the file or to standard output.

10026 **Command Search and Execution**

10027 This description requires that the shell can execute shell scripts directly, even if the underlying
 10028 system does not support the common "#!" interpreter convention. That is, if file **foo** contains
 10029 shell commands and is executable, the following executes **foo**:

```

10030     ./foo

```

10031 The command search shown here does not match all historical implementations. A more typical
 10032 sequence has been:

- 10033 • Any built-in (special or regular)
- 10034 • Functions
- 10035 • Path search for executable files

10036 But there are problems with this sequence. Since the programmer has no idea in advance which
 10037 utilities might have been built into the shell, a function cannot be used to override portably a
 10038 utility of the same name. (For example, a function named *cd* cannot be written for many
 10039 historical systems.) Furthermore, the *PATH* variable is partially ineffective in this case, and only
 10040 a pathname with a slash can be used to ensure a specific executable file is invoked.

10041 After the *execve()* failure described, the shell normally executes the file as a shell script. Some
 10042 implementations, however, attempt to detect whether the file is actually a script and not an
 10043 executable from some other architecture. The method used by the KornShell is allowed by the
 10044 text that indicates non-text files may be bypassed.

10045 The sequence selected for the Shell and Utilities volume of IEEE Std 1003.1-200x acknowledges
 10046 that special built-ins cannot be overridden, but gives the programmer full control over which
 10047 versions of other utilities are executed. It provides a means of suppressing function lookup (via
 10048 the *command* utility) for the user's own functions and ensures that any regular built-ins or
 10049 functions provided by the implementation are under the control of the path search. The
 10050 mechanisms for associating built-ins or functions with executable files in the path are not
 10051 specified by the Shell and Utilities volume of IEEE Std 1003.1-200x, but the wording requires that
 10052 if either is implemented, the application is not able to distinguish a function or built-in from an
 10053 executable (other than in terms of performance, presumably). The implementation ensures that
 10054 all effects specified by the Shell and Utilities volume of IEEE Std 1003.1-200x resulting from the
 10055 invocation of the regular built-in or function (interaction with the environment, variables, traps,
 10056 and so on) are identical to those resulting from the invocation of an executable file.

10057 Examples

10058 Consider three versions of the *ls* utility:

- 10059 1. The application includes a shell function named *ls*.
- 10060 2. The user writes a utility named *ls* and puts it in **/fred/bin**.
- 10061 3. The example implementation provides *ls* as a regular shell built-in that is invoked (either
 10062 by the shell or directly by *exec*) when the path search reaches the directory **/posix/bin**.

10063 If *PATH*=**/posix/bin**, various invocations yield different versions of *ls*:

10064

10065

10066

10067

10068

10069

10070

Invocation	Version of <i>ls</i>
<i>ls</i> (from within application script)	(1) function
<i>command ls</i> (from within application script)	(3) built-in
<i>ls</i> (from within makefile called by application)	(3) built-in
<i>system("ls")</i>	(3) built-in
<i>PATH="/fred/bin:\$PATH" ls</i>	(2) user's version

10071 C.2.9.2 Pipelines

10072 Because pipeline assignment of standard input or standard output or both takes place before
 10073 redirection, it can be modified by redirection. For example:

```
10074 $ command1 2>&1 | command2
```

10075 sends both the standard output and standard error of *command1* to the standard input of
 10076 *command2*.

10077 The reserved word **!** allows more flexible testing using AND and OR lists.

10078 It was suggested that it would be better to return a non-zero value if any command in the
 10079 pipeline terminates with non-zero status (perhaps the bitwise-inclusive OR of all return values).
 10080 However, the choice of the last-specified command semantics are historical practice and would
 10081 cause applications to break if changed. An example of historical behavior:

```
10082 $ sleep 5 | (exit 4)
```

```
10083 $ echo $?
```

```
10084 4
```

```
10085 $ (exit 4) | sleep 5
```

```
10086 $ echo $?
```

```
10087 0
```

10088 C.2.9.3 Lists

10089 The equal precedence of "&&" and "||" is historical practice. The standard developers
 10090 evaluated the model used more frequently in high-level programming languages, such as C, to
 10091 allow the shell logical operators to be used for complex expressions in an unambiguous way, but
 10092 they could not allow historical scripts to break in the subtle way unequal precedence might
 10093 cause. Some arguments were posed concerning the "{" or "(" groupings that are required
 10094 historically. There are some disadvantages to these groupings:

- 10095 • The "(" can be expensive, as they spawn other processes on some implementations. This
 10096 performance concern is primarily an implementation issue.
- 10097 • The "{" braces are not operators (they are reserved words) and require a trailing space
 10098 after each '{', and a semicolon before each '}'. Most programmers (and certainly
 10099 interactive users) have avoided braces as grouping constructs because of the problematic
 10100 syntax required. Braces were not changed to operators because that would generate
 10101 compatibility issues even greater than the precedence question; braces appear outside the
 10102 context of a keyword in many shell scripts.

10103 IEEE PASC Interpretation 1003.2 #204 is applied, clarifying that the operators "&&" and "||"
 10104 are evaluated with left associativity.

10105 **Asynchronous Lists**

10106 The grammar treats a construct such as:

```
10107     foo & bar & bam &
```

10108 as one "asynchronous list", but since the status of each element is tracked by the shell, the term
 10109 "element of an asynchronous list" was introduced to identify just one of the **foo**, **bar**, or **bam**
 10110 portions of the overall list.

10111 Unless the implementation has an internal limit, such as {CHILD_MAX}, on the retained process
 10112 IDs, it would require unbounded memory for the following example:

```
10113     while true
10114     do       foo & echo $!
10115     done
```

10116 The treatment of the signals SIGINT and SIGQUIT with asynchronous lists is described in the
 10117 Shell and Utilities volume of IEEE Std 1003.1-200x, Section 2.11, Signals and Error Handling.

10118 Since the connection of the input to the equivalent of /dev/null is considered to occur before
 10119 redirections, the following script would produce no output:

```
10120     exec < /etc/passwd
10121     cat <&0 &
10122     wait
```

10123 **Sequential Lists**

10124 There is no additional rationale provided for this section.

- 10125 **AND Lists**
- 10126 There is no additional rationale provided for this section.
- 10127 **OR Lists**
- 10128 There is no additional rationale provided for this section.
- 10129 *C.2.9.4 Compound Commands*
- 10130 **Grouping Commands**
- 10131 The semicolon shown in *{compound-list;}* is an example of a control operator delimiting the } reserved word. Other delimiters are possible, as shown in the Shell and Utilities volume of IEEE Std 1003.1-200x, Section 2.10, Shell Grammar; <newline> is frequently used.
- 10132
- 10133
- 10134 A proposal was made to use the <do-done> construct in all cases where command grouping in the current process environment is performed, identifying it as a construct for the grouping commands, as well as for shell functions. This was not included because the shell already has a grouping construct for this purpose ("{}"), and changing it would have been counter-productive.
- 10135
- 10136
- 10137
- 10138
- 10139 **For Loop**
- 10140 The format is shown with generous usage of <newline>s. See the grammar in the Shell and Utilities volume of IEEE Std 1003.1-200x, Section 2.10, Shell Grammar for a precise description of where <newline>s and semicolons can be interchanged.
- 10141
- 10142
- 10143 Some historical implementations support ' { ' and ' } ' as substitutes for **do** and **done**. The standard developers chose to omit them, even as an obsolescent feature. (Note that these substitutes were only for the **for** command; the **while** and **until** commands could not use them historically because they are followed by compound-lists that may contain "{ . . . }" grouping commands themselves.)
- 10144
- 10145
- 10146
- 10147
- 10148 The reserved word pair **do** . . . **done** was selected rather than **do** . . . **od** (which would have matched the spirit of **if** . . . **fi** and **case** . . . **esac**) because *od* is already the name of a standard utility.
- 10149
- 10150
- 10151 PASC Interpretation 1003.2 #169 has been applied changing the grammar.
- 10152 **Case Conditional Construct**
- 10153 An optional left parenthesis before *pattern* was added to allow numerous historical KornShell scripts to conform. At one time, using the leading parenthesis was required if the **case** statement was to be embedded within a "\$ ()" command substitution; this is no longer the case with the POSIX shell. Nevertheless, many historical scripts use the left parenthesis, if only because it makes matching-parenthesis searching easier in *vi* and other editors. This is a relatively simple implementation change that is upward-compatible for all scripts.
- 10154
- 10155
- 10156
- 10157
- 10158
- 10159 Consideration was given to requiring *break* inside the *compound-list* to prevent falling through to the next pattern action list. This was rejected as being nonexistent practice. An interesting undocumented feature of the KornShell is that using ";&" instead of ";;" as a terminator causes the exact opposite behavior—the flow of control continues with the next *compound-list*.
- 10160
- 10161
- 10162
- 10163 The pattern ' * ', given as the last pattern in a **case** construct, is equivalent to the default case in a C-language **switch** statement.
- 10164

10165 The grammar shows that reserved words can be used as patterns, even if one is the first word on
10166 a line. Obviously, the reserved word **esac** cannot be used in this manner.

10167 **If Conditional Construct**

10168 The precise format for the command syntax is described in the Shell and Utilities volume of
10169 IEEE Std 1003.1-200x, Section 2.10, Shell Grammar.

10170 **While Loop**

10171 The precise format for the command syntax is described in the Shell and Utilities volume of
10172 IEEE Std 1003.1-200x, Section 2.10, Shell Grammar.

10173 **Until Loop**

10174 The precise format for the command syntax is described in the Shell and Utilities volume of
10175 IEEE Std 1003.1-200x, Section 2.10, Shell Grammar.

10176 *C.2.9.5 Function Definition Command*

10177 The description of functions in an early proposal was based on the notion that functions should
10178 behave like miniature shell scripts; that is, except for sharing variables, most elements of an
10179 execution environment should behave as if they were a new execution environment, and
10180 changes to these should be local to the function. For example, traps and options should be reset
10181 on entry to the function, and any changes to them do not affect the traps or options of the caller.
10182 There were numerous objections to this basic idea, and the opponents asserted that functions
10183 were intended to be a convenient mechanism for grouping common commands that were to be
10184 executed in the current execution environment, similar to the execution of the *dot* special built-
10185 in.

10186 It was also pointed out that the functions described in that early proposal did not provide a local
10187 scope for everything a new shell script would, such as the current working directory, or *umask*,
10188 but instead provided a local scope for only a few select properties. The basic argument was that
10189 if a local scope is needed for the execution environment, the mechanism already existed: the
10190 application can put the commands in a new shell script and call that script. All historical shells
10191 that implemented functions, other than the KornShell, have implemented functions that operate
10192 in the current execution environment. Because of this, traps and options have a global scope
10193 within a shell script. Local variables within a function were considered and included in another
10194 early proposal (controlled by the special built-in *local*), but were removed because they do not fit
10195 the simple model developed for functions and because there was some opposition to adding yet
10196 another new special built-in that was not part of historical practice. Implementations should
10197 reserve the identifier *local* (as well as *typeset*, as used in the KornShell) in case this local variable
10198 mechanism is adopted in a future version of IEEE Std 1003.1-200x.

10199 A separate issue from the execution environment of a function is the availability of that function
10200 to child shells. A few objectors maintained that just as a variable can be shared with child shells
10201 by exporting it, so should a function. In early proposals, the *export* command therefore had a *-f*
10202 flag for exporting functions. Functions that were exported were to be put into the environment
10203 as *name()=value* pairs, and upon invocation, the shell would scan the environment for these and
10204 automatically define these functions. This facility was strongly opposed and was omitted. Some
10205 of the arguments against exportable functions were as follows:

- 10206 • There was little historical practice. The Ninth Edition shell provided them, but there was
10207 controversy over how well it worked.

10208 • There are numerous security problems associated with functions appearing in the
 10209 environment of a user and overriding standard utilities or the utilities owned by the
 10210 application.

10211 • There was controversy over requiring *make* to import functions, where it has historically used
 10212 an *exec* function for many of its command line executions.

10213 • Functions can be big and the environment is of a limited size. (The counter-argument was
 10214 that functions are no different from variables in terms of size: there can be big ones, and there
 10215 can be small ones—and just as one does not export huge variables, one does not export huge
 10216 functions. However, this might not apply to the average shell-function writer, who typically
 10217 writes much larger functions than variables.)

10218 As far as can be determined, the functions in the Shell and Utilities volume of
 10219 IEEE Std 1003.1-200x match those in System V. Earlier versions of the KornShell had two
 10220 methods of defining functions:

```
10221       function fname { compound-list }
```

10222 and:

```
10223       fname() { compound-list }
```

10224 The latter used the same definition as the Shell and Utilities volume of IEEE Std 1003.1-200x, but
 10225 differed in semantics, as described previously. The current edition of the KornShell aligns the
 10226 latter syntax with the Shell and Utilities volume of IEEE Std 1003.1-200x and keeps the former as
 10227 is.

10228 The name space for functions is limited to that of a *name* because of historical practice.
 10229 Complications in defining the syntactic rules for the function definition command and in dealing
 10230 with known extensions such as the "@()" usage in the KornShell prevented the name space
 10231 from being widened to a *word*. Using functions to support synonyms such as the "!!" and '%'
 10232 usage in the C shell is thus disallowed to conforming applications, but acceptable as an
 10233 extension. For interactive users, the aliasing facilities in the Shell and Utilities volume of
 10234 IEEE Std 1003.1-200x should be adequate for this purpose. It is recognized that the name space
 10235 for utilities in the file system is wider than that currently supported for functions, if the portable
 10236 filename character set guidelines are ignored, but it did not seem useful to mandate extensions
 10237 in systems for so little benefit to conforming applications.

10238 The "()" in the function definition command consists of two operators. Therefore, intermixing
 10239 <blank>s with the *fname*, '(', and ')' is allowed, but unnecessary.

10240 An example of how a function definition can be used wherever a simple command is allowed:

```
10241       # If variable i is equal to "yes",  

  10242       # define function foo to be ls -l  

  10243       #  

  10244       [ "$i" = yes ] && foo() {  

  10245           ls -l  

  10246       }
```

10247 **C.2.10 Shell Grammar**

10248 There are several subtle aspects of this grammar where conventional usage implies rules about
10249 the grammar that in fact are not true.

10250 For *compound_list*, only the forms that end in a *separator* allow a reserved word to be recognized,
10251 so usually only a *separator* can be used where a compound list precedes a reserved word (such as
10252 **Then, Else, Do** and **Rbrace**). Explicitly requiring a separator would disallow such valid (if rare)
10253 statements as:

```
10254     if (false) then (echo x) else (echo y) fi
```

10255 See the Note under special grammar rule 1.

10256 Concerning the third sentence of rule (1) (“Also, if the parser ...”):

10257 • This sentence applies rather narrowly: when a compound list is terminated by some clear
10258 delimiter (such as the closing **fi** of an inner **if_clause**) then it would apply; where the
10259 compound list might continue (as in after a ‘;’), rule (7a) (and consequently the first
10260 sentence of rule (1)) would apply. In many instances the two conditions are identical, but this
10261 part of rule (1) does not give license to treating a **WORD** as a reserved word unless it is in a
10262 place where a reserved word has to appear.

10263 • The statement is equivalent to requiring that when the LR(1) lookahead set contains exactly
10264 one reserved word, it must be recognized if it is present. (Here “LR(1)” refers to the
10265 theoretical concepts, not to any real parser generator.)

10266 For example, in the construct below, and when the parser is at the point marked with ‘^’,
10267 the only next legal token is **then** (this follows directly from the grammar rules):

```
10268     if if...fi then ... fi
10269         ^
```

10270 At that point, the **then** must be recognized as a reserved word.

10271 (Depending on the parser generator actually used, “extra” reserved words may be in some
10272 lookahead sets. It does not really matter if they are recognized, or even if any possible
10273 reserved word is recognized in that state, because if it is recognized and is not in the
10274 (theoretical) LR(1) lookahead set, an error is ultimately detected. In the example above, if
10275 some other reserved word (for example, **while**) is also recognized, an error occurs later.

10276 This is approximately equivalent to saying that reserved words are recognized after other
10277 reserved words (because it is after a reserved word that this condition occurs), but avoids the
10278 “except for ...” list that would be required for **case**, **for**, and so on. (Reserved words are of
10279 course recognized anywhere a *simple_command* can appear, as well. Other rules take care of
10280 the special cases of non-recognition, such as rule (4) for **case** statements.)

10281 Note that the body of here-documents are handled by token recognition (see the Shell and
10282 Utilities volume of IEEE Std 1003.1-200x, Section 2.3, Token Recognition) and do not appear in
10283 the grammar directly. (However, the here-document I/O redirection operator is handled as part
10284 of the grammar.)

10285 The start symbol of the grammar (**complete_command**) represents either input from the
10286 command line or a shell script. It is repeatedly applied by the interpreter to its input and
10287 represents a single “chunk” of that input as seen by the interpreter.

10288 *C.2.10.1 Shell Grammar Lexical Conventions*

10289 There is no additional rationale provided for this section.

10290 *C.2.10.2 Shell Grammar Rules*

10291 There is no additional rationale provided for this section.

10292 **C.2.11 Signals and Error Handling**

10293 There is no additional rationale provided for this section.

10294 **C.2.12 Shell Execution Environment**10295 Some implementations have implemented the last stage of a pipeline in the current environment |
10296 so that commands such as: |10297 `command | read foo` |10298 set variable **foo** in the current environment. This extension is allowed, but not required;
10299 therefore, a shell programmer should consider a pipeline to be in a subshell environment, but
10300 not depend on it.10301 In early proposals, the description of execution environment failed to mention that each
10302 command in a multiple command pipeline could be in a subshell execution environment. For
10303 compatibility with some historical shells, the wording was phrased to allow an implementation
10304 to place any or all commands of a pipeline in the current environment. However, this means that
10305 a POSIX application must assume each command is in a subshell environment, but not depend
10306 on it.10307 The wording about shell scripts is meant to convey the fact that describing “trap actions” can
10308 only be understood in the context of the shell command language. Outside of this context, such
10309 as in a C-language program, signals are the operative condition, not traps.10310 **C.2.13 Pattern Matching Notation**10311 Pattern matching is a simpler concept and has a simpler syntax than REs, as the former is
10312 generally used for the manipulation of filenames, which are relatively simple collections of
10313 characters, while the latter is generally used to manipulate arbitrary text strings of potentially
10314 greater complexity. However, some of the basic concepts are the same, so this section points
10315 liberally to the detailed descriptions in the Base Definitions volume of IEEE Std 1003.1-200x,
10316 Chapter 9, Regular Expressions.10317 *C.2.13.1 Patterns Matching a Single Character*10318 Both quoting and escaping are described here because pattern matching must work in three
10319 separate circumstances:

- 10320 1. Calling directly upon the shell, such as in pathname expansion or in a
- case**
- statement. All
-
- 10321 of the following match the string or file
- abc**
- :

10322 `abc "abc" a"b" c a\bc a[b]c a["b"]c a[\b]c a["\b"]c a?c a*c` |

10323 The following do not:

10324 `"a?c" a*c a\[b]c` |

- 10325 2. Calling a utility or function without going through a shell, as described for
- find*
- and the
-
- 10326
- fnmatch()*
- function defined in the System Interfaces volume of IEEE Std 1003.1-200x.

10327 3. Calling utilities such as *find*, *cpio*, *tar*, or *pax* through the shell command line. In this case,
10328 shell quote removal is performed before the utility sees the argument. For example, in:

```
10329 find /bin -name "e\c[\h]o" -print
```

10330 after quote removal, the backslashes are presented to *find* and it treats them as escape
10331 characters. Both precede ordinary characters, so the *c* and *h* represent themselves and *echo*
10332 would be found on many historical systems (that have it in */bin*). To find a filename that
10333 contained shell special characters or pattern characters, both quoting and escaping are
10334 required, such as:

```
10335 pax -r ... "*a\(\?"
```

10336 to extract a filename ending with "a(?".

10337 Conforming applications are required to quote or escape the shell special characters (sometimes
10338 called metacharacters). If used without this protection, syntax errors can result or
10339 implementation extensions can be triggered. For example, the KornShell supports a series of
10340 extensions based on parentheses in patterns.

10341 The restriction on a circumflex in a bracket expression is to allow implementations that support
10342 pattern matching using the circumflex as the negation character in addition to the exclamation
10343 mark. A conforming application must use something like "[\^!]" to match either character.

10344 C.2.13.2 Patterns Matching Multiple Characters

10345 Since each asterisk matches zero or more occurrences, the patterns "a*b" and "a**b" have
10346 identical functionality.

10347 Examples

10348 a[bc] Matches the strings "ab" and "ac".

10349 a*d Matches the strings "ad", "abd", and "abcd", but not the string "abc".

10350 a*d* Matches the strings "ad", "abcd", "abcdef", "aaaad", and "adddd".

10351 *a*d Matches the strings "ad", "abcd", "efabcd", "aaaad", and "adddd".

10352 C.2.13.3 Patterns Used for Filename Expansion

10353 The caveat about a slash within a bracket expression is derived from historical practice. The
10354 pattern "a[b/c]d" does not match such pathnames as **abd** or **a/d**. On some implementations
10355 (including those conforming to the Single UNIX Specification), it matched a pathname of
10356 literally "a[b/c]d". On other systems, it produced an undefined condition (an unescaped '['
10357 used outside a bracket expression). In this version, the XSI behavior is now required.

10358 Filenames beginning with a period historically have been specially protected from view on
10359 UNIX systems. A proposal to allow an explicit period in a bracket expression to match a leading
10360 period was considered; it is allowed as an implementation extension, but a conforming
10361 application cannot make use of it. If this extension becomes popular in the future, it will be
10362 considered for a future version of the Shell and Utilities volume of IEEE Std 1003.1-200x.

10363 Historical systems have varied in their permissions requirements. To match **f*/bar** has required
10364 read permissions on the **f*** directories in the System V shell, but the Shell and Utilities volume of
10365 IEEE Std 1003.1-200x, the C shell, and KornShell require only search permissions.

10366 C.2.14 Special Built-In Utilities

10367 See the RATIONALE sections on the individual reference pages.

10368 C.3 Batch Environment Services and Utilities**10369 Scope of the Batch Environment Option**

10370 This section summarizes the deliberations of the IEEE P1003.15 (Batch Environment) working
10371 group in the development of the Batch Environment option, which covers a set of services and
10372 utilities defining a batch processing system.

10373 This informative section contains historical information concerning the contents of the
10374 amendment and describes why features were included or discarded by the working group.

10375 History of Batch Systems

10376 The supercomputing technical committee began as a “Birds Of a Feather” (BOF) at the January
10377 1987 Usenix meeting. There was enough general interest to form a supercomputing attachment
10378 to the /usr/group working groups. Several subgroups rapidly formed. Of those subgroups, the
10379 batch group was the most ambitious. The first early meetings were spent evaluating user needs
10380 and existing batch implementations.

10381 To evaluate user needs, individuals from the supercomputing community came and presented
10382 their needs. Common requests were flexibility, interoperability, control of resources, and ease-
10383 of-use. Backwards-compatibility was not an issue. The working group then evaluated some
10384 existing systems. The following different systems were evaluated:

- 10385 • PROD
- 10386 • Convex Distributed Batch
- 10387 • NQS
- 10388 • CTSS
- 10389 • MDQS from Ballistics Research Laboratory (BRL)

10390 Finally, NQS was chosen as a model because it satisfied not only the most user requirements, but
10391 because it was public domain, already implemented on a variety of hardware platforms, and
10392 networked-based.

10393 Historical Implementations of Batch Systems

10394 Deferred processing of work under the control of a scheduler has been a feature of most
10395 proprietary operating systems from the earliest days of multi-user systems in order to maximize
10396 utilization of the computer.

10397 The arrival of UNIX systems proved to be a dilemma to many hardware providers and users
10398 because it did not include the sophisticated batch facilities offered by the proprietary systems.
10399 This omission was rectified in 1986 by NASA Ames Research Center who developed the
10400 Network Queuing System (NQS) as a portable UNIX application that allowed the routing and
10401 processing of batch “jobs” in a network. To encourage its usage, the product was later put into
10402 the public domain. It was promptly picked up by UNIX hardware providers, and ported and
10403 developed for their respective hardware and UNIX implementations.

10404 Many major vendors, who traditionally offer a batch-dominated environment, ported the
 10405 public-domain product to their systems, customized it to support the capabilities of their
 10406 systems, and added many customer-requested features.

10407 Due to the strong hardware provider and customer acceptance of NQS, it was decided to use
 10408 NQS as the basis for the POSIX Batch Environment amendment in 1987. Other batch systems
 10409 considered at the time included CTSS, MDQS (a forerunner of NQS from the Ballistics Research
 10410 Laboratory), and PROD (a Los Alamos Labs development). None were thought to have both the
 10411 functionality and acceptability of NQS.

10412 **NQS Differences from the *at* utility**

10413 The base standard *at* and *batch* utilities are not sufficient to meet the batch processing needs in a
 10414 supercomputing environment and additional functionality in the areas of resource management,
 10415 job scheduling, system management, and control of output is required.

10416 **Batch Environment Option Definitions**

10417 The concept of a batch job is closely related to a session with a session leader. The main
 10418 difference is that a batch job does not have a controlling terminal. There has been much debate
 10419 over whether to use the term *request* or *job*. *Job* was the final choice because of the historical use
 10420 of this term in the batch environment.

10421 The current definition for job identifiers is not sufficient with the model of destinations. The
 10422 current definition is:

```
10423     sequence_number.originating_host
```

10424 Using the model of destination, a host may include multiple batch nodes, the location of which is
 10425 identified uniquely by a name or directory service. If the current definition is used, batch nodes
 10426 running on the same host would have to coordinate their use of sequence numbers, as sequence
 10427 numbers are assigned by the originating host. The alternative is to use the originating batch node
 10428 name instead of the originating host name.

10429 The reasons for wishing to run more than one batch system per host could be the following:

10430 A test and production batch system are maintained on a single host. This is most likely in a
 10431 development facility, but could also arise when a site is moving from one version to another.
 10432 The new batch system could be installed as a test version that is completely separate from the
 10433 production batch system, so that problems can be isolated to the test system. Requiring the batch
 10434 nodes to coordinate their use of sequence numbers creates a dependency between the two
 10435 nodes, and that defeats the purpose of running two nodes.

10436 A site has multiple departments using a single host, with different management policies. An
 10437 example of contention might be in job selection algorithms. One group might want a FIFO type
 10438 of selection, while another group wishes to use a more complex algorithm based on resource
 10439 availability. Again, requiring the batch nodes to coordinate is an unnecessary binding.

10440 The proposal eventually accepted was to replace originating host with originating batch node.
 10441 This supplies sufficient granularity to ensure unique job identifiers. If more than one batch node
 10442 is on a particular host, they each have their own unique name.

10443 The queue portion of a destination is not part of the job identifier as these are not required to be
 10444 unique between batch nodes. For instance, two batch nodes may both have queues called small,
 10445 medium, and large. It is only the batch node name that is uniquely identifiable throughout the
 10446 batch system. The queue name has no additional function in this context.

10447 Assume there are three batch nodes, each of which has its own name server. On batch node one,
10448 there are no queues. On batch node two, there are fifty queues. On batch node three, there are
10449 forty queues. The system administrator for batch node one does not have to configure queues,
10450 because there are none implemented. However, if a user wishes to send a job to either batch
10451 node two or three, the system administrator for batch node one must configure a destination
10452 that maps to the appropriate batch node and queue. If every queue is to be made accessible from
10453 batch node one, the system administrator has to configure ninety destinations.

10454 To avoid requiring this, there should be a mechanism to allow a user to separate the destination
10455 into a batch node name and a queue name. Then, an implementation that is configured to get to
10456 all the batch nodes does not need any more configuration to allow a user to get to all of the
10457 queues on all of the batch nodes. The node name is used to locate the batch node, while the
10458 queue name is sent unchanged to that batch node.

10459 The following are requirements that a destination identifier must be capable of providing:

- 10460 • The ability to direct a job to a queue in a particular batch node.
- 10461 • The ability to direct a job to a particular batch node.
- 10462 • The ability to group at a higher level than just one queue. This includes grouping similar
10463 queues across multiple batch nodes (this is a pipe queue today).
- 10464 • The ability to group batch nodes. This allows a user to submit a job to a group name with no
10465 knowledge of the batch node configuration. This also provides aliasing as a special case.
10466 Aliasing is a group containing only one batch node name. The group name is the alias.

10467 In addition, the administrator has the following requirements:

- 10468 • The ability to control access to the queues.
- 10469 • The ability to control access to the batch nodes.
- 10470 • The ability to control access to groups of queues (pipe queues).
- 10471 • The ability to configure retry time intervals and durations.

10472 The requirements of the user are met by destination as explained in the following:

10473 The user has the ability to specify a queue name, which is known only to the batch node
10474 specified. There is no configuration of these queues required on the submitting node.

10475 The user has the ability to specify a batch node whose name is network-unique. The
10476 configuration required is that the batch node be defined as an application, just as other
10477 applications such as FTP are configured.

10478 Once a job reaches a queue, it can again become a user of the batch system. The batch node can
10479 choose to send the job to another batch node or queue or both. In other words, the routing is at
10480 an application level, and it is up to the batch system to choose where the job will be sent.
10481 Configuration is up to the batch node where the queue resides. This provides grouping of
10482 queues across batch nodes or within a batch node. The user submits the job to a queue, which by
10483 definition routes the job to other queues or nodes or both.

10484 A node name may be given to a naming service, which returns multiple addresses as opposed to
10485 just one. This provides grouping at a batch node level. This is a local issue, meaning that the
10486 batch node must choose only one of these addresses. The list of addresses is not sent with the
10487 job, and once the job is accepted on another node, there is no connection between the list and the
10488 job. The requirements of the administrator are met by destination as explained in the following:

10489 The control of queues is a batch system issue, and will be done using the batch administrative
10490 utilities.

- 10491 The control of nodes is a network issue, and will be done through whatever network facilities
10492 are available.
- 10493 The control of access to groups of queues (pipe queues) is covered by the control of any other
10494 queue. The fact that the job may then be sent to another destination is not relevant.
- 10495 The propagation of a job across more than one point-to-point connection was dropped because
10496 of its complexity and because all of the issues arising from this capability could not be resolved.
10497 It could be provided as additional functionality at some time in the future.
- 10498 The addition of *network* as a defined term was done to clarify the difference between a network
10499 of batch nodes as opposed to a network of hosts. A network of batch nodes is referred to as a
10500 batch system. The network refers to the actual host configuration. A single host may have
10501 multiple batch nodes.
- 10502 In the absence of a standard network naming convention, this option establishes its own
10503 convention for the sake of consistency and expediency. This is subject to change, should a future
10504 working group develop a standard naming convention for network pathnames.

10505 C.3.1 Batch General Concepts

- 10506 During the development of the Batch Environment option, a number of topics were discussed at
10507 length which influenced the wording of the normative text but could not be included in the final
10508 text. The following items are some of the most significant terms and concepts of those discussed:
- 10509 • Small and Consistent Command Set

10510 Often, conventional utilities from UNIX systems have a very complicated utility syntax and
10511 usage. This can often result in confusion and errors when trying to use them. The Batch
10512 Environment option utility set, on the other hand, has been paired to a small set of robust
10513 utilities with an orthogonal calling sequence.
 - 10514 • Checkpoint/Restart

10515 This feature permits an already executing process to checkpoint or save its contents. Some
10516 implementations permit this at both the batch utility level; for example, checkpointing this
10517 job upon its abnormal termination or from within the job itself via a system call. Support of
10518 checkpoint/restart is optional. A conscious, careful effort was made to make the *qsub* and
10519 *qmgr* utilities consistently refer to checkpoint/restart as optional functionality.
 - 10520 • Rerunability

10521 When a user submits a job for batch processing, they can designate it “rerunnable” in that it
10522 will automatically resume execution from the start of the job if the machine on which it was
10523 executing crashes for some reason. The decision on whether the job will be rerun or not is
10524 entirely up to the submitter of the job and no decisions will be made within the batch system.
10525 A job that is rerunnable and has been submitted with the proper checkpoint/restart switch
10526 will first be checkpointed and execution begun from that point. Furthermore, use of the
10527 implementation-defined checkpoint/restart feature will be not be defined in this context.
 - 10528 • Error Codes

10529 All utilities exit with error status zero (0) if successful, one (1) if a user error occurred, and
10530 two (2) for an internal Batch Environment option error.
 - 10531 • Level of Portability

10532 Portability is specified at both the user, operator, and administrator levels. A conforming
10533 batch implementation prevents identical functionality and behavior at all these levels.
10534 Additionally, portable batch shell scripts with embedded Batch Environment option utilities

- 10535 adds an additional level of portability.
- 10536
- Resource Specification
- 10537 A small set of globally understood resources, such as memory and CPU time, is specified. All
10538 conforming batch implementations are able to process them in a manner consistent with the
10539 yet-to-be-developed resource management model. Resources not in this amendment set are
10540 ignored and passed along as part of the argument stream of the utility.
- 10541
- Queue Position
- 10542 Queue position is the place a job occupies in a queue. It is dependent on a variety of factors
10543 such as submission time and priority. Since priority may be affected by the implementation
10544 of fair share scheduling, the definition of queue position is implementation-defined.
- 10545
- Queue ID
- 10546 A numerical queue ID is an external requirement for purposes of accounting. The
10547 identification number was chosen over queue name for processing convenience.
- 10548
- Job ID
- 10549 A common notion of “jobs” is a collection of processes whose process group cannot be
10550 altered and is used for resource management and accounting. This concept is
10551 implementation-defined and, as such, has been omitted from the batch amendment.
- 10552
- Bytes *versus* Words
- 10553 Except for one case, bytes are used as the standard unit for memory size. Furthermore, the
10554 definition of a word varies from machine to machine. Therefore, bytes will be the default unit
10555 of memory size.
- 10556
- Regular Expressions
- 10557 The standard definition of regular expressions is much too broad to be used in the batch
10558 utility syntax. All that is needed is a simple concept of “all”; for example, delete all my jobs
10559 from the named queue. For this reason, regular expressions have been eliminated from the
10560 batch amendment.
- 10561
- Display Privacy
- 10562 How much data should be displayed locally through functions? Local policy dictates the
10563 amount of privacy. Library functions must be used to create and enforce local policy.
10564 Network and local *qstats* must reflect the policy of the server machine.
- 10565
- Remote Host Naming Convention
- 10566 It was decided that host names would be a maximum of 255 characters in length, with at
10567 most 15 characters being shown in displays. The 255 character limit was chosen because it is
10568 consistent with BSD. The 15-character limit was an arbitrary decision.
- 10569
- Network Administration
- 10570 Network administration is important, but is outside the scope of the batch amendment.
10571 Network administration could done with *rsh*. However, authentication becomes two-sided.
- 10572
- Network Administration Philosophy
- 10573 Keep it simple. Centralized management should be possible. For example, Los Alamos needs
10574 a dumb set of CPUs to be managed by a central system *versus* several independently-
10575 managed systems as is the general case for the Batch Environment option.

- 10576 • Operator Utility Defaults (that is, Default Host, User, Account, and so on)
- 10577 It was decided that usability would override orthogonality and syntactic consistency.
- 10578 • The Batch System Manager and Operator Distinction
- 10579 The distinction between manager and operator is that operators can only control the flow of
- 10580 jobs. A manager can alter the batch system configuration in addition to job flow. POSIX
- 10581 makes a distinction between user and system administrator but goes no further. The
- 10582 concepts of manager and operator privileges fall under local policy. The distinction between
- 10583 manager and operator is historical in batch environments, and the Batch Environment option
- 10584 has continued that distinction.
- 10585 • The Batch System Administrator
- 10586 An administrator is equivalent to a batch system manager.

10587 C.3.2 Batch Services

- 10588 This rationale is provided as informative rather than normative text, to avoid placing
 10589 requirements on implementors regarding the use of symbolic constants, but at the same time to
 10590 give implementors a preferred practice for assigning values to these constants to promote
 10591 interoperability.
- 10592 The *Checkpoint* and *Minimum_Cpu_Interval* attributes induce a variety of behavior depending
 10593 upon their values. Some jobs cannot or should not be checkpointed. Other users will simply
 10594 need to ensure job continuation across planned downtimes; for example, scheduled preventive
 10595 maintenance. For users consuming expensive resources, or for jobs that run longer than the
 10596 mean time between failures, however, periodic checkpointing may be essential. However,
 10597 system administrators must be able to set minimum checkpoint intervals on a queue-by-queue
 10598 basis to guard against; for example, naive users specifying interval values too small on memory
 10599 intensive jobs. Otherwise, system overhead would adversely affect performance.
- 10600 The use of symbolic constants, such as `NO_CHECKPOINT`, was introduced to lend a degree of
 10601 formalism and portability to this option.
- 10602 Support for checkpointing is optional for servers. However, clients must provide for the `-c`
 10603 option, since in a distributed environment the job may run on a server that does provide such
 10604 support, even if the host of the client does not support the checkpoint feature.
- 10605 If the user does not specify the `-c` option, the default action is left unspecified by this option.
 10606 Some implementations may wish to do checkpointing by default; others may wish to checkpoint
 10607 only under an explicit request from the user.
- 10608 The *Priority* attribute has been made non-optional. All clients already had been required to
 10609 support the `-p` option. The concept of prioritization is common in historical implementations.
 10610 The default priority is left to the server to establish.
- 10611 The *Hold_Types* attribute has been modified to allow for implementation-defined hold types to
 10612 be passed to a batch server.
- 10613 It was the intent of the IEEE P1003.15 working group to mandate the support for the
 10614 *Resource_List* attribute in this option by referring to another amendment, specifically P1003.1a.
 10615 However, during the development of P1003.1a this was excluded. As such this requirement has
 10616 been removed from the normative text.
- 10617 The *Shell_Path* attribute has been modified to accept a list of shell paths that are associated with
 10618 a host. The name of the attribute has been changed to *Shell_Path_List*.

10619 **C.3.3 Common Behavior for Batch Environment Utilities**

10620 This section was defined to meet the goal of a “Small and Consistent Command Set” for this
10621 option.

10622 **C.4 Utilities**

10623 For the utilities included in IEEE Std 1003.1-200x, see the RATIONALE sections on the individual
10624 reference pages.

10625 **Exclusion of Utilities**

10626 The set of utilities contained in IEEE Std 1003.1-200x is drawn from the base documents, with
10627 one addition: the *c99* utility. This section contains rationale for some of the deliberations that led
10628 to this set of utilities, and why certain utilities were excluded.

10629 Many utilities were evaluated by the standard developers; more historical utilities were
10630 excluded from the base documents than included. The following list contains many common
10631 UNIX system utilities that were not included as mandatory utilities, in the UPE, in the XSI
10632 extension, or in one of the software development groups. It is logistically difficult for this
10633 rationale to distribute correctly the reasons for not including a utility among the various utility
10634 options. Therefore, this section covers the reasons for all utilities not included in
10635 IEEE Std 1003.1-200x.

10636 This rationale is limited to a discussion of only those utilities actively or indirectly evaluated by
10637 the standard developers of the base documents, rather than the list of all known UNIX utilities
10638 from all its variants.

10639 *adb* The intent of the various software development utilities was to assist in the
10640 installation (rather than the actual development and debugging) of applications.
10641 This utility is primarily a debugging tool. Furthermore, many useful aspects of *adb*
10642 are very hardware-specific.

10643 *as* Assemblers are hardware-specific and are included implicitly as part of the
10644 compilers in IEEE Std 1003.1-200x.

10645 *banner* The only known use of this command is as part of the *lp* printer header pages. It
10646 was decided that the format of the header is implementation-defined, so this utility
10647 is superfluous to application portability.

10648 *calendar* This reminder service program is not useful to conforming applications. |

10649 *cancel* The *lp* (line printer spooling) system specified is the most basic possible and did
10650 not need this level of application control.

10651 *chroot* This is primarily of administrative use, requiring superuser privileges.

10652 *col* No utilities defined in IEEE Std 1003.1-200x produce output requiring such a filter.
10653 The *nroff* text formatter is present on many historical systems and will continue to
10654 remain as an extension; *col* is expected to be shipped by all the systems that ship
10655 *nroff*.

10656 *cpio* This has been replaced by *pax*, for reasons explained in the rationale for that utility.

10657 *cpp* This is subsumed by *c99*.

10658 *cu* This utility is terminal-oriented and is not useful from shell scripts or typical
10659 application programs. |

10660	<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.
10661		
10662		
10663		
10664	<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.
10665		
10666		
10667	<i>dis</i>	Disassemblers are hardware-specific.
10668	<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> .
10669		
10670		
10671		
10672		
10673		
10674		
10675		
10676	<i>ld</i>	This is subsumed by <i>c99</i> .
10677	<i>line</i>	The functionality of <i>line</i> can be provided with <i>read</i> .
10678	<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.
10679		
10680		
10681		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.
10682		
10683		
10684		
10685		
10686		
10687		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).
10688		
10689		
10690		
10691		
10692		
10693		
10694	<i>login</i>	This utility is terminal-oriented and is not useful from shell scripts or typical application programs.
10695		
10696	<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.
10697		
10698	<i>lpstat</i>	The <i>lp</i> system specified is the most basic possible and did not need this level of application control.
10699		
10700	<i>mail</i>	This utility was omitted in favor of <i>mailx</i> because there was a considerable functionality overlap between the two.
10701		
10702	<i>mknod</i>	This was omitted in favor of <i>mkfifo</i> , as <i>mknod</i> has too many implementation-defined functions.
10703		

10704	<i>news</i>	This utility is terminal-oriented and is not useful from shell scripts or typical application programs.
10705		
10706	<i>pack</i>	This compression program was considered inferior to <i>compress</i> .
10707	<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:
10708		
10709		<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.
10710		
10711		
10712		<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.
10713		
10714		
10715		
10716		
10717		
10718		<ul style="list-style-type: none"> • This is really a subject for a system administration working group or a security working group.
10719		
10720	<i>pcat</i>	This compression program was considered inferior to <i>zcat</i> .
10721	<i>pg</i>	This duplicated many of the features of the <i>more</i> pager, which was preferred by the standard developers.
10722		
10723	<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.
10724		
10725		
10726	RCS	RCS was originally considered as part of a version control utilities portion of the scope. However, this aspect was abandoned by the standard developers. SCCS is now included as an optional part of the XSI extension.
10727		
10728		
10729	<i>red</i>	Restricted editor. This was not considered by the standard developers because it never provided the level of security restriction required.
10730		
10731	<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.
10732		
10733		
10734	<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.
10735		
10736		
10737		
10738	<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.
10739		
10740		
10741	<i>shar</i>	Any of the numerous “shell archivers” were excluded because they did not meet the requirement of existing practice.
10742		
10743	<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.
10744		
10745		
10746	<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.
10747		

10748		This utility is primarily a debugging tool.
10749	<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.
10750		
10751		
10752		
10753	<i>su</i>	This utility is not useful from shell scripts or typical application programs. (There was also sentiment to avoid security-related utilities.)
10754		
10755	<i>sum</i>	This utility was renamed <i>cksum</i> .
10756	<i>tar</i>	This has been replaced by <i>pax</i> , for reasons explained in the rationale for that utility.
10757	<i>tsort</i>	This utility is an aid in creating an implementation-defined detail of object libraries that the standard developers did not feel required standardization.
10758		
10759	<i>unpack</i>	This compression program was considered inferior to <i>uncompress</i> .
10760	<i>wall</i>	This utility is terminal-oriented and is not useful from shell scripts or typical application programs. It is generally used only by system administrators.
10761		

10763 / *Rationale (Informative)*

10764 **Part D:**

10765 **Portability Considerations**

10766 *The Open Group*

10767 *The Institute of Electrical and Electronics Engineers, Inc.*

Portability Considerations (Informative)

10769

10770 This section contains information to satisfy various international requirements:

- 10771 • Section D.1 describes perceived user requirements.
- 10772 • Section D.2 (on page 3558) indicates how the facilities of IEEE Std 1003.1-200x satisfy those
10773 requirements.
- 10774 • Section D.3 (on page 3565) offers guidance to writers of profiles on how the configurable
10775 options, limits, and optional behavior of IEEE Std 1003.1-200x should be cited in profiles.

10776 D.1 User Requirements

10777 This section describes the user requirements that were perceived by the developers of
10778 IEEE Std 1003.1-200x. The primary source for these requirements was an analysis of historical
10779 practice in widespread use, as typified by the base documents listed in Section A.1.1 (on page
10780 3293).

10781 IEEE Std 1003.1-200x addresses the needs of users requiring open systems solutions for source
10782 code portability of applications. It currently addresses users requiring open systems solutions
10783 for source-code portability of applications involving multi-programming and process
10784 management (creating processes, signaling, and so on); access to files and directories in a
10785 hierarchy of file systems (opening, reading, writing, deleting files, and so on); access to
10786 asynchronous communications ports and other special devices; access to information about
10787 other users of the system; facilities supporting applications requiring bounded (realtime)
10788 response.

10789 The following users are identified for IEEE Std 1003.1-200x:

- 10790 • Those employing applications written in high-level languages, such as C, Ada, or FORTRAN.
- 10791 • Users who desire conforming applications that do not necessarily require the characteristics
10792 of high-level languages (for example, the speed of execution of compiled languages or the
10793 relative security of source code intellectual property inherent in the compilation process).
- 10794 • Users who desire conforming applications that can be developed quickly and can be
10795 modified readily without the use of compilers and other system components that may be
10796 unavailable on small systems or those without special application development capabilities.
- 10797 • Users who interact with a system to achieve general-purpose time-sharing capabilities
10798 common to most business or government offices or academic environments: editing, filing,
10799 inter-user communications, printing, and so on.
- 10800 • Users who develop applications for POSIX-conformant systems.
- 10801 • Users who develop applications for UNIX systems.

10802 An acknowledged restriction on applicable users is that they are limited to the group of
10803 individuals who are familiar with the style of interaction characteristic of historically-derived
10804 systems based on one of the UNIX operating systems (as opposed to other historical systems
10805 with different models, such as MS/DOS, Macintosh, VMS, MVS, and so on). Typical users
10806 would include program developers, engineers, or general-purpose time-sharing users.

10807 The requirements of users of IEEE Std 1003.1-200x can be summarized as a single goal:
10808 *application source portability*. The requirements of the user are stated in terms of the requirements
10809 of portability of applications. This in turn becomes a requirement for a standardized set of
10810 syntax and semantics for operations commonly found on many operating systems.

10811 The following sections list the perceived requirements for application portability.

10812 **D.1.1 Configuration Interrogation**

10813 An application must be able to determine whether and how certain optional features are
10814 provided and to identify the system upon which it is running, so that it may appropriately adapt
10815 to its environment.

10816 Applications must have sufficient information to adapt to varying behaviors of the system.

10817 **D.1.2 Process Management**

10818 An application must be able to manage itself, either as a single process or as multiple processes.
10819 Applications must be able to manage other processes when appropriate.

10820 Applications must be able to identify, control, create, and delete processes, and there must be
10821 communication of information between processes and to and from the system.

10822 Applications must be able to use multiple flows of control with a process (threads) and
10823 synchronize operations between these flows of control.

10824 **D.1.3 Access to Data**

10825 Applications must be able to operate on the data stored on the system, access it, and transmit it
10826 to other applications. Information must have protection from unauthorized or accidental access
10827 or modification.

10828 **D.1.4 Access to the Environment**

10829 Applications must be able to access the external environment to communicate their input and
10830 results.

10831 **D.1.5 Access to Determinism and Performance Enhancements**

10832 Applications must have sufficient control of resource allocation to ensure the timeliness of
10833 interactions with external objects.

10834 **D.1.6 Operating System-Dependent Profile**

10835 The capabilities of the operating system may make certain optional characteristics of the base
10836 language in effect no longer optional, and this should be specified.

10837 D.1.7 I/O Interaction

10838 The interaction between the C language I/O subsystem (*stdio*) and the I/O subsystem of
10839 IEEE Std 1003.1-200x must be specified.

10840 D.1.8 Internationalization Interaction

10841 The effects of the environment of IEEE Std 1003.1-200x on the internationalization facilities of the
10842 C language must be specified.

10843 D.1.9 C-Language Extensions

10844 Certain functions in the C language must be extended to support the additional capabilities
10845 provided by IEEE Std 1003.1-200x.

10846 D.1.10 Command Language

10847 Users should be able to define procedures that combine simple tools and/or applications into
10848 higher-level components that perform to the specific needs of the user. The user should be able
10849 to store, recall, use, and modify these procedures. These procedures should employ a powerful
10850 command language that is used for recurring tasks in conforming applications (scripts) in the
10851 same way that it is used interactively to accomplish one-time tasks. The language and the
10852 utilities that it uses must be consistent between systems to reduce errors and retraining.

10853 D.1.11 Interactive Facilities

10854 Use the system to accomplish individual tasks at an interactive terminal. The interface should be
10855 consistent, intuitive, and offer usability enhancements to increase the productivity of terminal
10856 users, reduce errors, and minimize retraining costs. Online documentation or usage assistance
10857 should be available.

10858 D.1.12 Accomplish Multiple Tasks Simultaneously

10859 Access applications and interactive facilities from a single terminal without requiring serial
10860 execution: switch between multiple interactive tasks; schedule one-time or periodic background
10861 work; display the status of all work in progress or scheduled; influence the priority scheduling of
10862 work, when authorized.

10863 D.1.13 Complex Data Manipulation

10864 Manipulate data in files in complex ways: sort, merge, compare, translate, edit, format, pattern
10865 match, select subsets (strings, columns, fields, rows, and so on). These facilities should be
10866 available to both conforming applications and interactive users.

10867 D.1.14 File Hierarchy Manipulation

10868 Create, delete, move/rename, copy, backup/archive, and display files and directories. These
10869 facilities should be available to both conforming applications and interactive users.

10870 D.1.15 Locale Configuration

10871 Customize applications and interactive sessions for the cultural and language conventions of the
10872 user. Employ a wide variety of standard character encodings. These facilities should be available
10873 to both conforming applications and interactive users. |

10874 D.1.16 Inter-User Communication

10875 Send messages or transfer files to other users on the same system or other systems on a network.
10876 These facilities should be available to both conforming applications and interactive users. |

10877 D.1.17 System Environment

10878 Display information about the status of the system (activities of users and their interactive and
10879 background work, file system utilization, system time, configuration, and presence of optional
10880 facilities) and the environment of the user (terminal characteristics, and so on). Inform the
10881 system operator/administrator of problems. Control access to user files and other resources.

10882 D.1.18 Printing

10883 Output files on a variety of output device classes, accessing devices on local or network-
10884 connected systems. Control (or influence) the formatting, priority scheduling, and output
10885 distribution of work. These facilities should be available to both conforming applications and |
10886 interactive users.

10887 D.1.19 Software Development

10888 Develop (create and manage source files, compile/interpret, debug) portable open systems
10889 applications and package them for distribution to, and updating of, other systems.

10890 D.2 Portability Capabilities

10891 This section describes the significant portability capabilities of IEEE Std 1003.1-200x and
10892 indicates how the user requirements listed in Section D.1 (on page 3555) are addressed. The
10893 capabilities are listed in the same format as the preceding user requirements; they are
10894 summarized below:

- 10895 • Configuration Interrogation
- 10896 • Process Management
- 10897 • Access to Data
- 10898 • Access to the Environment
- 10899 • Access to Determinism and Performance Enhancements
- 10900 • Operating System-Dependent Profile
- 10901 • I/O Interaction
- 10902 • Internationalization Interaction
- 10903 • C-Language Extensions
- 10904 • Command Language
- 10905 • Interactive Facilities

- 10906 • Accomplish Multiple Tasks Simultaneously
- 10907 • Complex Data Manipulation
- 10908 • File Hierarchy Manipulation
- 10909 • Locale Configuration
- 10910 • Inter-User Communication
- 10911 • System Environment
- 10912 • Printing
- 10913 • Software Development

10914 **D.2.1 Configuration Interrogation**

10915 The *uname()* operation provides basic identification of the system. The *sysconf()*, *pathconf()*, and
 10916 *fpathconf()* functions and the *getconf* utility provide means to interrogate the implementation to
 10917 determine how to adapt to the environment in which it is running. These values can be either
 10918 static (indicating that all instances of the implementation have the same value) or dynamic
 10919 (indicating that different instances of the implementation have the different values, or that the
 10920 value may vary for other reasons, such as reconfiguration).

10921 **Unsatisfied Requirements**

10922 None directly. However, as new areas are added, there will be a need for additional capability in
 10923 this area.

10924 **D.2.2 Process Management**

10925 The *fork()*, *exec* family, and *spawn()* functions provide for the creation of new processes or the
 10926 insertion of new applications into existing processes. The *_Exit()*, *_exit()*, *exit()*, and *abort()*
 10927 functions allow for the termination of a process by itself. The *wait()* and *waitpid()* functions
 10928 allow one process to deal with the termination of another.

10929 The *times()* function allows for basic measurement of times used by a process. Various
 10930 functions, including *fstat()*, *getegid()*, *geteuid()*, *getgid()*, *getrgid()*, *getgrnam()*, *getlogin()*,
 10931 *getpid()*, *getppid()*, *getpwnam()*, *getpwuid()*, *getuid()*, *lstat()*, and *stat()*, provide for access to the
 10932 identifiers of processes and the identifiers and names of owners of processes (and files).

10933 The various functions operating on environment variables provide for communication of
 10934 information (primarily user-configurable defaults) from a parent to child processes.

10935 The operations on the current working directory control and interrogate the directory from
 10936 which relative filename searches start. The *umask()* function controls the default protections
 10937 applied to files created by the process.

10938 The *alarm()*, *pause()*, *sleep()*, *ualarm()*, and *usleep()* operations allow the process to suspend until
 10939 a timer has expired or to be notified when a period of time has elapsed. The *time()* operation
 10940 interrogates the current time and date.

10941 The signal mechanism provides for communication of events either from other processes or
 10942 from the environment to the application, and the means for the application to control the effect
 10943 of these events. The mechanism provides for external termination of a process and for a process
 10944 to suspend until an event occurs. The mechanism also provides for a value to be associated with
 10945 an event.

10946 Job control provides a means to group processes and control them as groups, and to control their
 10947 access to the function between the user and the system (the *controlling terminal*). It also provides
 10948 the means to suspend and resume processes.

10949 The Process Scheduling option provides control of the scheduling and priority of a process.

10950 The Message Passing option provides a means for interprocess communication involving small
 10951 amounts of data.

10952 The Memory Management facilities provide control of memory resources and for the sharing of
 10953 memory. This functionality is mandatory on XSI-conformant systems.

10954 The Threads facilities provide multiple flows of control with a process (threads),
 10955 synchronization between threads, association of data with threads, and controlled cancelation of
 10956 threads.

10957 The XSI interprocess communications functionality provide an alternate set of facilities to
 10958 manipulate semaphores, message queues, and shared memory. These are provided on XSI-
 10959 conformant systems to support conforming applications developed to run on UNIX systems.

10960 D.2.3 Access to Data

10961 The *open()*, *close()*, *fclose()*, *fopen()*, and *pipe()* functions provide for access to files and data.
 10962 Such files may be regular files, interprocess data channels (pipes), or devices. Additional types
 10963 of objects in the file system are permitted and are being contemplated for standardization.

10964 The *access()*, *chmod()*, *chown()*, *dup()*, *dup2()*, *fchmod()*, *fcntl()*, *fstat()*, *ftruncate()*, *lstat()*,
 10965 *readlink()*, *realpath()*, *stat()*, and *utime()* functions allow for control and interrogation of file and
 10966 file-related objects, (including symbolic links) and their ownership, protections, and timestamps.

10967 The *fgetc()*, *fputc()*, *fread()*, *fseek()*, *fsetpos()*, *fwrite()*, *getc()*, *getch()*, *lseek()*, *putchar()*, *putc()*,
 10968 *read()*, and *write()* functions provide for data transfer from the application to files (in all their
 10969 forms).

10970 The *closedir()*, *link()*, *mkdir()*, *opendir()*, *readdir()*, *rename()*, *rmdir()*, *rewinddir()*, and *unlink()*
 10971 functions provide for a complete set of operations on directories. Directories can arbitrarily
 10972 contain other directories, and a single file can be mentioned in more than one directory.

10973 The file-locking mechanism provides for advisory locking (protection during transactions) of
 10974 ranges of bytes (in effect, records) in a file.

10975 The *confstr()*, *fpathconf()*, *pathconf()*, and *sysconf()* functions provide for enquiry as to the
 10976 behavior of the system where variability is permitted.

10977 The Synchronized Input and Output option provides for assured commitment of data to media.

10978 The Asynchronous Input and Output option provides for initiation and control of asynchronous
 10979 data transfers.

10980 D.2.4 Access to the Environment

10981 The operations and types in the Base Definitions volume of IEEE Std 1003.1-200x, Chapter 11,
 10982 General Terminal Interface are provided for access to asynchronous serial devices. The primary
 10983 intended use for these is the controlling terminal for the application (the interaction point
 10984 between the user and the system). They are general enough to be used to control any
 10985 asynchronous serial device. The functions are also general enough to be used with many other
 10986 device types as a user interface when some emulation is provided.

10987 Less detailed access is provided for other device types, but in many instances an application
10988 need not know whether an object in the file system is a device or a regular file to operate
10989 correctly.

10990 **Unsatisfied Requirements**

10991 Detailed control of common device classes, specifically magnetic tape, is not provided.

10992 **D.2.5 Bounded (Realtime) Response**

10993 The Realtime Signals Extension provides queued signals and the prioritization of the handling of
10994 signals. The SCHED_FIFO, SCHED_SPORADIC, and SCHED_RR scheduling policies provide
10995 control over processor allocation. The Semaphores option provides high-performance
10996 synchronization. The Memory Management functions provide memory locking for control of
10997 memory allocation, file mapping for high-performance, and shared memory for high-
10998 performance interprocess communication. The Message Passing option provides for interprocess
10999 communication without being dependent on shared memory.

11000 The Timers option provides a high resolution function called *nanosleep()* with a finer resolution
11001 than the *sleep()* function.

11002 The Typed Memory Objects option, the Monotonic Clock option, and the Timeouts option
11003 provide further facilities for applications to use to obtain predictable bounded response.

11004 **D.2.6 Operating System-Dependent Profile**

11005 IEEE Std 1003.1-200x makes no distinction between text and binary files. The values of
11006 EXIT_SUCCESS and EXIT_FAILURE are further defined.

11007 **Unsatisfied Requirements**

11008 None known, but the ISO C standard may contain some additional options that could be
11009 specified.

11010 **D.2.7 I/O Interaction**

11011 IEEE Std 1003.1-200x defines how each of the ISO C standard *stdio* functions interact with the
11012 POSIX.1 operations, typically specifying the behavior in terms of POSIX.1 operations.

11013 **Unsatisfied Requirements**

11014 None.

11015 **D.2.8 Internationalization Interaction**

11016 The IEEE Std 1003.1-200x environment operations provide a means to define the environment
11017 for *setlocale()* and time functions such as *ctime()*. The *tzset()* function is provided to set time
11018 conversion information.

11019 The *nl_langinfo()* function is provided as an XSI extension to query locale-specific cultural
11020 settings.

11021 **Unsatisfied Requirements**

11022 None.

11023 **D.2.9 C-Language Extensions**11024 The *setjmp()* and *longjmp()* functions are not defined to be cognizant of the signal masks defined
11025 for POSIX.1. The *sigsetjmp()* and *siglongjmp()* functions are provided to fill this gap.11026 The *_setjmp()* and *_longjmp()* functions are provided as XSI extensions to support historic
11027 practice.11028 **Unsatisfied Requirements**

11029 None.

11030 **D.2.10 Command Language**11031 The shell command language, as described in Shell and Utilities volume of IEEE Std 1003.1-200x,
11032 Chapter 2, Shell Command Language, is a common language useful in batch scripts, through an
11033 API to high-level languages (for the C-Language Binding option, *system()* and *popen()*) and
11034 through an interactive terminal (see the *sh* utility). The shell language has many of the
11035 characteristics of a high-level language, but it has been designed to be more suitable for user
11036 terminal entry and includes interactive debugging facilities. Through the use of pipelining,
11037 many complex commands can be constructed from combinations of data filters and other
11038 common components. Shell scripts can be created, stored, recalled, and modified by the user
11039 with simple editors.11040 In addition to the basic shell language, the following utilities offer features that simplify and
11041 enhance programmatic access to the utilities and provide features normally found only in high-
11042 level languages: *basename*, *bc*, *command*, *dirname*, *echo*, *env*, *expr*, *false*, *printf*, *read*, *sleep*, *tee*, *test*,
11043 *time**,² *true*, *wait*, *xargs*, and all of the special built-in utilities in the Shell and Utilities volume of
11044 IEEE Std 1003.1-200x, Section 2.14, Special Built-In Utilities .11045 **Unsatisfied Requirements**

11046 None.

11047 **D.2.11 Interactive Facilities**11048 The utilities offer a common style of command-line interface through conformance to the Utility
11049 Syntax Guidelines (see the Base Definitions volume of IEEE Std 1003.1-200x, Section 12.2, Utility
11050 Syntax Guidelines) and the common utility defaults (see the Shell and Utilities volume of
11051 IEEE Std 1003.1-200x, Section 1.11, Utility Description Defaults). The *sh* utility offers an
11052 interactive command-line history and editing facility. The following utilities in the User
11053 Portability Utilities option have been customized for interactive use: *alias*, *ex*, *fc*, *mailx*, *more*, *talk*,
11054 *vi*, *unalias*, and *write*; the *man* utility offers online access to system documentation.

11055 _____

11056 2. The utilities listed with an asterisk here and later in this section are present only on systems which support the User Portability
11057 Utilities option. There may be further restrictions on the utilities offered with various configuration option combinations; see the
11058 individual utility descriptions. |

11059 **Unsatisfied Requirements**

11060 The command line interface to individual utilities is as intuitive and consistent as historical
 11061 practice allows. Work underway based on graphical user interfaces may be more suitable for
 11062 novice or occasional users of the system.

11063 **D.2.12 Accomplish Multiple Tasks Simultaneously**

11064 The shell command language offers background processing through the asynchronous list
 11065 command form; see the Shell and Utilities volume of IEEE Std 1003.1-200x, Section 2.9, Shell
 11066 Commands. The *nohup* utility makes background processing more robust and usable. The *kill*
 11067 utility can terminate background jobs. When the User Portability Utilities option is supported,
 11068 the following utilities allow manipulation of jobs: *bg*, *fg*, and *jobs*. Also, if the User Portability
 11069 Utilities option is supported, the following can support periodic job scheduling, control, and
 11070 display: *at*, *batch*, *crontab*, *nice*, *ps*, and *renice*.

11071 **Unsatisfied Requirements**

11072 Terminals with multiple windows may be more suitable for some multi-tasking interactive uses
 11073 than the job control approach in IEEE Std 1003.1-200x. See the comments on graphical user
 11074 interfaces in Section D.2.11 (on page 3562). The *nice* and *renice* utilities do not necessarily take
 11075 advantage of complex system scheduling algorithms that are supported by the realtime options
 11076 within IEEE Std 1003.1-200x.

11077 **D.2.13 Complex Data Manipulation**

11078 The following utilities address user requirements in this area: *asa*, *awk*, *bc*, *cmp*, *comm*, *csplit**, *cut*,
 11079 *dd*, *diff*, *ed*, *ex**, *expand**, *expr*, *find*, *fold*, *grep*, *head*, *join*, *od*, *paste*, *pr*, *printf*, *sed*, *sort*, *split**, *tabs**, *tail*,
 11080 *tr*, *unexpand**, *uniq*, *uudecode**, *uuencode**, and *wc*.

11081 **Unsatisfied Requirements**

11082 Sophisticated text formatting utilities, such as *troff* or *TeX*, are not included. Standards work in
 11083 the area of SGML may satisfy this.

11084 **D.2.14 File Hierarchy Manipulation**

11085 The following utilities address user requirements in this area: *basename*, *cd*, *chgrp*, *chmod*, *chown*,
 11086 *cksum*, *cp*, *dd*, *df**, *diff*, *dirname*, *du**, *find*, *ls*, *ln*, *mkdir*, *mkfifo*, *mv*, *patch**, *pathchk*, *pax*, *pwd*, *rm*, *rmdir*,
 11087 *test*, and *touch*.

11088 **Unsatisfied Requirements**

11089 Some graphical user interfaces offer more intuitive file manager components that allow file
 11090 manipulation through the use of icons for novice users.

11091 D.2.15 Locale Configuration

11092 The standard utilities are affected by the various *LC_* variables to achieve locale-dependent
11093 operation: character classification, collation sequences, regular expressions and shell pattern
11094 matching, date and time formats, numeric formatting, and monetary formatting. When the
11095 POSIX2_LOCALEDEF option is supported, applications can provide their own locale definition
11096 files. The following utilities address user requirements in this area: *date*, *ed*, *ex**, *find*, *grep*, *locale*,
11097 *localedef*, *more**, *sed*, *sh*, *sort*, *tr*, *uniq*, and *vi**.

11098 The *iconv()*, *iconv_close()*, and *iconv_open()* functions are available to allow an application to
11099 convert character data between supported character sets.

11100 The *genccat* utility and the *catopen()*, *catclose()*, and *catgets()* functions for message catalog
11101 manipulation are available on XSI-conformant systems.

11102 Unsatisfied Requirements

11103 Some aspects of multi-byte character and state-encoded character encodings have not yet been
11104 addressed. The C-language functions, such as *getopt()*, are generally limited to single-byte
11105 characters. The effect of the *LC_MESSAGES* variable on message formats is only suggested at
11106 this time.

11107 D.2.16 Inter-User Communication

11108 The following utilities address user requirements in this area: *cksum*, *mailx**, *mesg**, *patch**, *pax*,
11109 *talk**, *uudecode**, *uuencode**, *who**, and *write**.

11110 The historical UUCP utilities are included on XSI-conformant systems.

11111 Unsatisfied Requirements

11112 None.

11113 D.2.17 System Environment

11114 The following utilities address user requirements in this area: *chgrp*, *chmod*, *chown*, *df**, *du**, *env*,
11115 *getconf*, *id*, *logger*, *logname*, *mesg**, *newgrp**, *ps**, *stty*, *tput**, *tty*, *umask*, *uname*, and *who**.

11116 The *closelog()*, *openlog()*, *setlogmask()*, and *syslog()* functions provide System Logging facilities
11117 on XSI-conformant systems; these are analogous to the *logger* utility.

11118 Unsatisfied Requirements

11119 None.

11120 D.2.18 Printing

11121 The following utilities address user requirements in this area: *pr* and *lp*.

11122 **Unsatisfied Requirements**

11123 There are no features to control the formatting or scheduling of the print jobs.

11124 **D.2.19 Software Development**

11125 The following utilities address user requirements in this area: *ar*, *asa*, *awk*, *c99*, *ctags**, *fort77*,
11126 *getconf*, *getopts*, *lex*, *localedef*, *make*, *nm**, *od*, *patch**, *pax*, *strings**, *strip*, *time**, and *yacc*.

11127 The *system()*, *popen()*, *pclose()*, *regcomp()*, *regexec()*, *regerror()*, *regfree()*, *fnmatch()*, *getopt()*,
11128 *glob()*, *globfree()*, *wordexp()*, and *wordfree()* functions allow C-language programmers to access
11129 some of the interfaces used by the utilities, such as argument processing, regular expressions,
11130 and pattern matching.

11131 The SCCS source-code control system utilities are available on systems supporting the XSI
11132 Development option.

11133 **Unsatisfied Requirements**

11134 There are no language-specific development tools related to languages other than C and
11135 FORTRAN. The C tools are more complete and varied than the FORTRAN tools. There is no
11136 data dictionary or other CASE-like development tools.

11137 **D.2.20 Future Growth**

11138 It is arguable whether or not all functionality to support applications is potentially within the
11139 scope of IEEE Std 1003.1-200x. As a simple matter of practicality, it cannot be. Areas such as |
11140 graphics, application domain-specific functionality, windowing, and so on, should be in unique |
11141 standards. As such, they are properly “Unsatisfied Requirements” in terms of providing fully |
11142 conforming applications, but ones which are outside the scope of IEEE Std 1003.1-200x. |

11143 However, as the standards evolve, certain functionality once considered “exotic” enough to be |
11144 part of a separate standard become common enough to be included in a core standard such as |
11145 this. Realtime and networking, for example, have both moved from separate standards (with |
11146 much difficult cross-referencing) into IEEE Std 1003.1 over time, and although no specific areas |
11147 have been identified for inclusion in future revisions, such inclusions seem likely. |

11148 **D.3 Profiling Considerations**

11149 This section offers guidance to writers of profiles on how the configurable options, limits, and
11150 optional behavior of IEEE Std 1003.1-200x should be cited in profiles. Profile writers should
11151 consult the general guidance in POSIX.0 when writing POSIX Standardized Profiles.

11152 The information in this section is an inclusive list of features that should be considered by profile |
11153 writers. Subsetting of IEEE Std 1003.1-200x should follow the Base Definitions volume of |
11154 IEEE Std 1003.1-200x, Section 2.1.5.1, Subprofiling Considerations. A set of profiling options is |
11155 described in Appendix E (on page 3579). |

11156 **D.3.1 Configuration Options**

11157 There are two set of options suggested by IEEE Std 1003.1-200x: those for POSIX-conforming
 11158 systems and those for X/Open System Interface (XSI) conformance. The requirements for XSI
 11159 conformance are documented in the Base Definitions volume of IEEE Std 1003.1-200x and not
 11160 discussed further here, as they superset the POSIX conformance requirements.

11161 **D.3.2 Configuration Options (Shell and Utilities)**

11162 There are three broad optional configurations for the Shell and Utilities volume of
 11163 IEEE Std 1003.1-200x: basic execution system, development system, and user portability
 11164 interactive system. The options to support these, and other minor configuration options, are
 11165 listed in the Base Definitions volume of IEEE Std 1003.1-200x, Chapter 2, Conformance. Profile
 11166 writers should consult the following list and the comments concerning user requirements
 11167 addressed by various components in Section D.2 (on page 3558).

11168 POSIX2_UPE

11169 The system supports the User Portability Utilities option.

11170 This option is a requirement for a user portability interactive system. It is required
 11171 frequently except for those systems, such as embedded realtime or dedicated application
 11172 systems, that support little or no interactive time-sharing work by users or operators. XSI-
 11173 conformant systems support this option.

11174 POSIX2_SW_DEV

11175 The system supports the Software Development Utilities option.

11176 This option is required by many systems, even those in which actual software development
 11177 does not occur. The *make* utility, in particular, is required by many application software
 11178 packages as they are installed onto the system. If POSIX2_C_DEV is supported,
 11179 POSIX2_SW_DEV is almost a mandatory requirement because of *ar* and *make*.

11180 POSIX2_C_BIND

11181 The system supports the C-Language Bindings option.

11182 This option is required on some implementations developing complex C applications or on
 11183 any system installing C applications in source form that require the functions in this option.
 11184 The *system()* and *popen()* functions, in particular, are widely used by applications; the
 11185 others are rather more specialized.

11186 POSIX2_C_DEV

11187 The system supports the C-Language Development Utilities option.

11188 This option is required by many systems, even those in which actual C-language software
 11189 development does not occur. The *c99* utility, in particular, is required by many application
 11190 software packages as they are installed onto the system. The *lex* and *yacc* utilities are used
 11191 less frequently.

11192 POSIX2_FORT_DEV

11193 The system supports the FORTRAN Development Utilities option

11194 As with C, this option is needed on any system developing or installing FORTRAN
 11195 applications in source form.

11196 POSIX2_FORT_RUN

11197 The system supports the FORTRAN Runtime Utilities option.

11198 This option is required for some FORTRAN applications that need the *asa* utility to convert
 11199 Hollerith printing statement output. It is unknown how frequently this occurs.

- 11200 POSIX2_LOCALEDEF
 11201 The system supports the creation of locales.
- 11202 This option is needed if applications require their own customized locale definitions to
 11203 operate. It is presently unknown whether many applications are dependent on this.
 11204 However, the option is virtually mandatory for systems in which internationalized
 11205 applications are developed.
- 11206 XSI-conformant systems support this option.
- 11207 POSIX2_PBS
 11208 The system supports the Batch Environment option.
- 11209 POSIX2_PBS_ACCOUNTING
 11210 The system supports the optional feature of accounting within the Batch Environment
 11211 option. It will be required in servers that implement the optional feature of accounting.
- 11212 POSIX2_PBS_CHECKPOINT
 11213 The systems supports the optional feature of checkpoint/restart within the Batch
 11214 Environment option.
- 11215 POSIX2_PBS_LOCATE
 11216 The system supports the optional feature of locating batch jobs within the Batch
 11217 Environment option.
- 11218 POSIX2_PBS_MESSAGE
 11219 The system supports the optional feature of sending messages to batch jobs within the
 11220 Batch Environment option.
- 11221 POSIX2_PBS_TRACK
 11222 The system supports the optional feature of tracking batch jobs within the Batch
 11223 Environment option.
- 11224 POSIX2_CHAR_TERM
 11225 The system supports at least one terminal type capable of all operations described in
 11226 IEEE Std 1003.1-200x.
- 11227 On systems with POSIX2_UPE, this option is almost always required. It was developed
 11228 solely to allow certain specialized vendors and user applications to bypass the requirement
 11229 for general-purpose asynchronous terminal support. For example, an application and
 11230 system that was suitable for block-mode terminals, such as IBM 3270s, would not need this
 11231 option.
- 11232 XSI-conformant systems support this option.

11233 D.3.3 Configurable Limits

- 11234 Very few of the limits need to be increased for profiles. No profile can cite lower values.
- 11235 {POSIX2_BC_BASE_MAX}
 11236 {POSIX2_BC_DIM_MAX}
 11237 {POSIX2_BC_SCALE_MAX}
 11238 {POSIX2_BC_STRING_MAX}
- 11239 No increase is anticipated for any of these *bc* values, except for very specialized applications
 11240 involving huge numbers.
- 11241 {POSIX2_COLL_WEIGHTS_MAX}
- 11242 Some natural languages with complex collation requirements require an increase from the
 11243 default 2 to 4; no higher numbers are anticipated.

- 11244 {POSIX2_EXPR_NEST_MAX}
 11245 No increase is anticipated.
- 11246 {POSIX2_LINE_MAX}
 11247 This number is much larger than most historical applications have been able to use. At some
 11248 future time, applications may be rewritten to take advantage of even larger values.
- 11249 {POSIX2_RE_DUP_MAX}
 11250 No increase is anticipated.
- 11251 {POSIX2_VERSION}
 11252 This is actually not a limit, but a standard version stamp. Generally, a profile should specify
 11253 Shell and Utilities volume of IEEE Std 1003.1-200x, Chapter 2, Shell Command Language by
 11254 name in the normative references section, not this value.

11255 **D.3.4 Configuration Options (System Interfaces)**

- 11256 {NGROUPS_MAX}
 11257 A non-zero value indicates that the implementation supports supplementary groups.
- 11258 This option is needed where there is a large amount of shared use of files, but where a
 11259 certain amount of protection is needed. Many profiles³ are known to require this option; it
 11260 should only be required if needed, but it should never be prohibited.
- 11261 _POSIX_ADVISORY_INFO
 11262 The system provides advisory information for file management.
- 11263 This option allows the application to specify advisory information that can be used to
 11264 achieve better or even deterministic response time in file manager or input and output
 11265 operations.
- 11266 _POSIX_ASYNCHRONOUS_IO
 11267 The system provides concurrent process execution and input and output transfers.
- 11268 This option was created to support historical systems that did not provide the feature. It
 11269 should only be required if needed, but it should never be prohibited.
- 11270 _POSIX_BARRIERS
 11271 The system supports barrier synchronization.
- 11272 This option was created to allow efficient synchronization of multiple parallel threads in
 11273 multi-processor systems in which the operation is supported in part by the hardware
 11274 architecture.
- 11275 _POSIX_CHOWN_RESTRICTED
 11276 The system restricts the right to “give away” files to other users.
- 11277 This option should be carefully investigated before it is required. Some applications expect
 11278 that they can change the ownership of files in this way. It is provided where either security
 11279 or system account requirements cause this ability to be a problem. It is also known to be
 11280 specified in many profiles.

11281 _____ |
 11282 3. There are no formally approved profiles of IEEE Std 1003.1-200x at the time of publication; the reference here is to various |
 11283 profiles generated by private bodies or governments. |

- 11284 `_POSIX_CLOCK_SELECTION`
11285 The system supports the Clock Selection option.
- 11286 This option allows applications to request a high resolution sleep in order to suspend a
11287 thread during a relative time interval, or until an absolute time value, using the desired
11288 clock. It also allows the application to select the clock used in a `pthread_cond_timedwait()`
11289 function call.
- 11290 `_POSIX_CPUTIME`
11291 The system supports the Process CPU-Time Clocks option.
- 11292 This option allows applications to use a new clock that measures the execution times of
11293 processes or threads, and the possibility to create timers based upon these clocks, for
11294 runtime detection (and treatment) of execution time overruns.
- 11295 `_POSIX_FSYNC`
11296 The system supports file synchronization requests.
- 11297 This option was created to support historical systems that did not provide the feature.
11298 Applications that are expecting guaranteed completion of their input and output operations
11299 should require the `_POSIX_SYNC_IO` option. This option should never be prohibited.
- 11300 XSI-conformant systems support this option.
- 11301 `_POSIX_IPV6`
11302 The system supports facilities related to Internet Protocol Version 6 (IPv6).
- 11303 This option was created to allow systems to transition to IPv6.
- 11304 `_POSIX_JOB_CONTROL`
11305 Job control facilities are mandatory in IEEE Std 1003.1-200x.
- 11306 The option was created primarily to support historical systems that did not provide the
11307 feature. Many existing profiles now require it; it should only be required if needed, but it
11308 should never be prohibited. Most applications that use it can run when it is not present,
11309 although with a degraded level of user convenience.
- 11310 `_POSIX_MAPPED_FILES`
11311 The system supports the mapping of regular files into the process address space.
- 11312 XSI-conformant systems support this option.
- 11313 Both this option and the Shared Memory Objects option provide shared access to memory
11314 objects in the process address space. The functions defined under this option provide the
11315 functionality of existing practice for mapping regular files. This functionality was deemed
11316 unnecessary, if not inappropriate, for embedded systems applications and, hence, is
11317 provided under this option. It should only be required if needed, but it should never be
11318 prohibited.
- 11319 `_POSIX_MEMLOCK`
11320 The system supports the locking of the address space.
- 11321 This option was created to support historical systems that did not provide the feature. It
11322 should only be required if needed, but it should never be prohibited.
- 11323 `_POSIX_MEMLOCK_RANGE`
11324 The system supports the locking of specific ranges of the address space.
- 11325 For applications that have well-defined sections that need to be locked and others that do
11326 not, IEEE Std 1003.1-200x supports an optional set of functions to lock or unlock a range of
11327 process addresses. The following are two reasons for having a means to lock down a

- 11328 specific range:
- 11329 1. An asynchronous event handler function that must respond to external events in a
11330 deterministic manner such that page faults cannot be tolerated
- 11331 2. An input/output “buffer” area that is the target for direct-to-process I/O, and the
11332 overhead of implicit locking and unlocking for each I/O call cannot be tolerated
- 11333 It should only be required if needed, but it should never be prohibited.
- 11334 _POSIX_MEMORY_PROTECTION
11335 The system supports memory protection.
- 11336 XSI-conformant systems support this option.
- 11337 The provision of this option typically imposes additional hardware requirements. It should
11338 never be prohibited.
- 11339 _POSIX_PRIORITIZED_IO
11340 The system provides prioritization for input and output operations.
- 11341 The use of this option may interfere with the ability of the system to optimize input and
11342 output throughput. It should only be required if needed, but it should never be prohibited.
- 11343 _POSIX_MESSAGE_PASSING
11344 The system supports the passing of messages between processes.
- 11345 This option was created to support historical systems that did not provide the feature. The
11346 functionality adds a high-performance XSI interprocess communication facility for local
11347 communication. It should only be required if needed, but it should never be prohibited.
- 11348 _POSIX_MONOTONIC_CLOCK
11349 The system supports the Monotonic Clock option.
- 11350 This option allows realtime applications to rely on a monotonically increasing clock that
11351 does not jump backwards, and whose value does not change except for the regular ticking
11352 of the clock.
- 11353 _POSIX_PRIORITY_SCHEDULING
11354 The system provides priority-based process scheduling.
- 11355 Support of this option provides predictable scheduling behavior, allowing applications to
11356 determine the order in which processes that are ready to run are granted access to a
11357 processor. It should only be required if needed, but it should never be prohibited.
- 11358 _POSIX_REALTIME_SIGNALS
11359 The system provides prioritized, queued signals with associated data values.
- 11360 This option was created to support historical systems that did not provide the features. It
11361 should only be required if needed, but it should never be prohibited.
- 11362 _POSIX_REGEX
11363 Support for regular expression facilities are mandatory in IEEE Std 1003.1-200x.
- 11364 _POSIX_SAVED_IDS
11365 Support for this feature is mandatory in IEEE Std 1003.1-200x.
- 11366 Certain classes of applications rely on it for proper operation, and there is no alternative
11367 short of giving the application root privileges on most implementations that did not provide
11368 _POSIX_SAVED_IDS.

- 11369 _POSIX_SEMAPHORES
11370 The system provides counting semaphores.
- 11371 This option was created to support historical systems that did not provide the feature. It
11372 should only be required if needed, but it should never be prohibited.
- 11373 _POSIX_SHARED_MEMORY_OBJECTS
11374 The system supports the mapping of shared memory objects into the process address space.
- 11375 Both this option and the Memory Mapped Files option provide shared access to memory
11376 objects in the process address space. The functions defined under this option provide the
11377 functionality of existing practice for shared memory objects. This functionality was deemed
11378 appropriate for embedded systems applications and, hence, is provided under this option. It
11379 should only be required if needed, but it should never be prohibited.
- 11380 _POSIX_SHELL
11381 Support for the *sh* utility command line interpreter is mandatory in IEEE Std 1003.1-200x.
- 11382 _POSIX_SPAWN
11383 The system supports the spawn option.
- 11384 This option provides applications with an efficient mechanism to spawn execution of a new
11385 process.
- 11386 _POSIX_SPINLOCKS
11387 The system supports spin locks.
- 11388 This option was created to support a simple and efficient synchronization mechanism for
11389 threads executing in multi-processor systems.
- 11390 _POSIX_SPORADIC_SERVER
11391 The system supports the sporadic server scheduling policy.
- 11392 This option provides applications with a new scheduling policy for scheduling aperiodic
11393 processes or threads in hard realtime applications.
- 11394 _POSIX_SYNCHRONIZED_IO
11395 The system supports guaranteed file synchronization.
- 11396 This option was created to support historical systems that did not provide the feature.
11397 Applications that are expecting guaranteed completion of their input and output operations
11398 should require this option, rather than the File Synchronization option. It should only be
11399 required if needed, but it should never be prohibited.
- 11400 _POSIX_THREADS
11401 The system supports multiple threads of control within a single process.
- 11402 This option was created to support historical systems that did not provide the feature.
11403 Applications written assuming a multi-threaded environment would be expected to require
11404 this option. It should only be required if needed, but it should never be prohibited.
- 11405 XSI-conformant systems support this option.
- 11406 _POSIX_THREAD_ATTR_STACKADDR
11407 The system supports specification of the stack address for a created thread.
- 11408 Applications may take advantage of support of this option for performance benefits, but
11409 dependence on this feature should be minimized. This option should never be prohibited.
- 11410 XSI-conformant systems support this option.

- 11411 _POSIX_THREAD_ATTR_STACKSIZE
11412 The system supports specification of the stack size for a created thread.
- 11413 Applications may require this option in order to ensure proper execution, but such usage
11414 limits portability and dependence on this feature should be minimized. It should only be
11415 required if needed, but it should never be prohibited.
- 11416 XSI-conformant systems support this option.
- 11417 _POSIX_THREAD_PRIORITY_SCHEDULING
11418 The system provides priority-based thread scheduling.
- 11419 Support of this option provides predictable scheduling behavior, allowing applications to
11420 determine the order in which threads that are ready to run are granted access to a processor.
11421 It should only be required if needed, but it should never be prohibited.
- 11422 _POSIX_THREAD_PRIO_INHERIT
11423 The system provides mutual exclusion operations with priority inheritance.
- 11424 Support of this option provides predictable scheduling behavior, allowing applications to
11425 determine the order in which threads that are ready to run are granted access to a processor.
11426 It should only be required if needed, but it should never be prohibited.
- 11427 _POSIX_THREAD_PRIO_PROTECT
11428 The system supports a priority ceiling emulation protocol for mutual exclusion operations.
- 11429 Support of this option provides predictable scheduling behavior, allowing applications to
11430 determine the order in which threads that are ready to run are granted access to a processor.
11431 It should only be required if needed, but it should never be prohibited.
- 11432 _POSIX_THREAD_PROCESS_SHARED
11433 The system provides shared access among multiple processes to synchronization objects.
- 11434 This option was created to support historical systems that did not provide the feature. It
11435 should only be required if needed, but it should never be prohibited.
- 11436 XSI-conformant systems support this option.
- 11437 _POSIX_THREAD_SAFE_FUNCTIONS
11438 The system provides thread-safe versions of all of the POSIX.1 functions.
- 11439 This option is required if the Threads option is supported. This is a separate option because
11440 thread-safe functions are useful in implementations providing other mechanisms for
11441 concurrency. It should only be required if needed, but it should never be prohibited.
- 11442 XSI-conformant systems support this option.
- 11443 _POSIX_THREAD_SPORADIC_SERVER
11444 The system supports the thread sporadic server scheduling policy.
- 11445 Support for this option provides applications with a new scheduling policy for scheduling
11446 aperiodic threads in hard realtime applications.
- 11447 _POSIX_TIMEOUTS
11448 The system provides timeouts for some blocking services.
- 11449 This option was created to provide a timeout capability to system services, thus allowing
11450 applications to include better error detection, and recovery capabilities.
- 11451 _POSIX_TIMERS
11452 The system provides higher resolution clocks with multiple timers per process.

11453 This option was created to support historical systems that did not provide the features. This
 11454 option is appropriate for applications requiring higher resolution timestamps or needing to
 11455 control the timing of multiple activities. It should only be required if needed, but it should
 11456 never be prohibited.

11457 `_POSIX_TRACE`

11458 The system supports the trace option.

11459 This option was created to allow applications to perform tracing.

11460 `_POSIX_TRACE_EVENT_FILTER`

11461 The system supports the trace event filter option.

11462 This option is dependent on support of the Trace option.

11463 `_POSIX_TRACE_INHERIT`

11464 The system supports the trace inherit option.

11465 This option is dependent on support of the Trace option.

11466 `_POSIX_TRACE_LOG`

11467 The system supports the trace log option.

11468 This option is dependent on support of the Trace option.

11469 `_POSIX_TYPED_MEMORY_OBJECTS`

11470 The system supports typed memory objects.

11471 This option was created to allow realtime applications to access different kinds of physical
 11472 memory, and allow processes in these applications to share portions of this memory.

11473 D.3.5 Configurable Limits

11474 In general, the configurable limits in the `<limits.h>` header defined in the Base Definitions
 11475 volume of IEEE Std 1003.1-200x have been set to minimal values; many applications or
 11476 implementations may require larger values. No profile can cite lower values.

11477 `{AIO_LISTIO_MAX}`

11478 The current minimum is likely to be inadequate for most applications. It is expected that
 11479 this value will be increased by profiles requiring support for list input and output
 11480 operations.

11481 `{AIO_MAX}`

11482 The current minimum is likely to be inadequate for most applications. It is expected that
 11483 this value will be increased by profiles requiring support for asynchronous input and
 11484 output operations.

11485 `{AIO_PRIO_DELTA_MAX}`

11486 The functionality associated with this limit is needed only by sophisticated applications. It
 11487 is not expected that this limit would need to be increased under a general-purpose profile.

11488 `{ARG_MAX}`

11489 The current minimum is likely to need to be increased for profiles, particularly as larger
 11490 amounts of information are passed through the environment. Many implementations are
 11491 believed to support larger values.

11492 `{CHILD_MAX}`

11493 The current minimum is suitable only for systems where a single user is not running
 11494 applications in parallel. It is significantly too low for any system also requiring windows,
 11495 and if `_POSIX_JOB_CONTROL` is specified, it should be raised.

- 11496 {CLOCKRES_MIN}
11497 It is expected that profiles will require a finer granularity clock, perhaps as fine as 1 μ s,
11498 represented by a value of 1 000 for this limit.
- 11499 {DELAYTIMER_MAX}
11500 It is believed that most implementations will provide larger values.
- 11501 {LINK_MAX}
11502 For most applications and usage, the current minimum is adequate. Many implementations
11503 have a much larger value, but this should not be used as a basis for raising the value unless
11504 the applications to be used require it.
- 11505 {LOGIN_NAME_MAX}
11506 This is not actually a limit, but an implementation parameter. No profile should impose a
11507 requirement on this value.
- 11508 {MAX_CANON}
11509 For most purposes, the current minimum is adequate. Unless high-speed burst serial
11510 devices are used, it should be left as is.
- 11511 {MAX_INPUT}
11512 See {MAX_CANON}.
- 11513 {MQ_OPEN_MAX}
11514 The current minimum should be adequate for most profiles.
- 11515 {MQ_PRIO_MAX}
11516 The current minimum corresponds to the required number of process scheduling priorities.
11517 Many realtime practitioners believe that the number of message priority levels ought to be
11518 the same as the number of execution scheduling priorities.
- 11519 {NAME_MAX}
11520 Many implementations now support larger values, and many applications and users
11521 assume that larger names can be used. Many existing profiles also specify a larger value.
11522 Specifying this value will reduce the number of conforming implementations, although this
11523 might not be a significant consideration over time. Values greater than 255 should not be
11524 required.
- 11525 {NGROUPS_MAX}
11526 The value selected will typically be 8 or larger.
- 11527 {OPEN_MAX}
11528 The historically common value for this has been 20. Many implementations support larger
11529 values. If applications that use larger values are anticipated, an appropriate value should be
11530 specified.
- 11531 {PAGESIZE}
11532 This is not actually a limit, but an implementation parameter. No profile should impose a
11533 requirement on this value.
- 11534 {PATH_MAX}
11535 Historically, the minimum has been either 1024 or indefinite, depending on the
11536 implementation. Few applications actually require values larger than 256, but some users
11537 may create file hierarchies that must be accessed with longer paths. This value should only
11538 be changed if there is a clear requirement.
- 11539 {PIPE_BUF}
11540 The current minimum is adequate for most applications. Historically, it has been larger. If
11541 applications that write single transactions larger than this are anticipated, it should be

11542 increased. Applications that write lines of text larger than this probably do not need it
11543 increased, as the text line is delimited by a newline.

11544 {POSIX_VERSION}
11545 This is actually not a limit, but a standard version stamp. Generally, a profile should specify
11546 IEEE Std 1003.1-200x by a name in the normative references section, not this value.

11547 {PTHREAD_DESTRUCTOR_ITERATIONS}
11548 It is unlikely that applications will need larger values to avoid loss of memory resources.

11549 {PTHREAD_KEYS_MAX}
11550 The current value should be adequate for most profiles.

11551 {PTHREAD_STACK_MIN}
11552 This should not be treated as an actual limit, but as an implementation parameter. No
11553 profile should impose a requirement on this value.

11554 {PTHREAD_THREADS_MAX}
11555 It is believed that most implementations will provide larger values.

11556 {RTSIG_MAX}
11557 The current limit was chosen so that the set of POSIX.1 signal numbers can fit within a 32-
11558 bit field. It is recognized that most existing implementations define many more signals than
11559 are specified in POSIX.1 and, in fact, many implementations have already exceeded 32
11560 signals (including the “null signal”). Support of {_POSIX_RTSIG_MAX} additional signals
11561 may push some implementations over the single 32-bit word line, but is unlikely to push
11562 any implementations that are already over that line beyond the 64 signal line.

11563 {SEM_NSEMS_MAX}
11564 The current value should be adequate for most profiles.

11565 {SEM_VALUE_MAX}
11566 The current value should be adequate for most profiles.

11567 {SSIZE_MAX}
11568 This limit reflects fundamental hardware characteristics (the size of an integer), and should
11569 not be specified unless it is clearly required. Extreme care should be taken to assure that
11570 any value that might be specified does not unnecessarily eliminate implementations
11571 because of accidents of hardware design.

11572 {STREAM_MAX}
11573 This limit is very closely related to {OPEN_MAX}. It should never be larger than
11574 {OPEN_MAX}, but could reasonably be smaller for application areas where most files are
11575 not accessed through *stdio*. Some implementations may limit {STREAM_MAX} to 20 but
11576 allow {OPEN_MAX} to be considerably larger. Such implementations should be allowed for
11577 if the applications permit.

11578 {TIMER_MAX}
11579 The current limit should be adequate for most profiles, but it may need to be larger for
11580 applications with a large number of asynchronous operations.

11581 {TTY_NAME_MAX}
11582 This is not actually a limit, but an implementation parameter. No profile should impose a
11583 requirement on this value.

11584 {TZNAME_MAX}
11585 The minimum has been historically adequate, but if longer timezone names are anticipated
11586 (particularly such values as UTC-1), this should be increased.

11587 **D.3.6 Optional Behavior**

11588 In IEEE Std 1003.1-200x, there are no instances of the terms unspecified, undefined,
11589 implementation-defined, or with the verbs “may” or “need not”, that the developers of
11590 IEEE Std 1003.1-200x anticipate or sanction as suitable for profile or test method citation. All of |
11591 these are merely warnings to conforming applications to avoid certain areas that can vary from |
11592 system to system, and even over time on the same system. In many cases, these terms are used
11593 explicitly to support extensions, but profiles should not anticipate and require such extensions; |
11594 future versions of IEEE Std 1003.1-200x may do so. |

11595	/ Rationale (Informative)	
11596	Part E:	
11597	Subprofiling Considerations	
11598	<i>The Open Group</i>	
11599	<i>The Institute of Electrical and Electronics Engineers, Inc.</i>	

Subprofiling Considerations (Informative)

11601

11602 This section contains further information to satisfy the requirement that the project scope enable
 11603 subprofiling of IEEE Std 1003.1-200x. The original intent was to have included a set of options
 11604 similar to the “Units of Functionality” contained in IEEE Std 1003.13-1998. However, as the
 11605 development of IEEE Std 1003.1-200x continued, the standard developers felt it premature to fix
 11606 these in normative text. The approach instead has been to include a general requirement in
 11607 normative text regarding subprofiling and to include an informative section (here) containing a
 11608 proposed set of subprofiling options.

11609 E.1 Subprofiling Option Groups

11610 The following Option Groups⁴ are defined to support profiling. Systems claiming support to
 11611 IEEE Std 1003.1-200x need not implement these options apart from the requirements stated in
 11612 the Base Definitions volume of IEEE Std 1003.1-200x, Section 2.1.3, POSIX Conformance. These
 11613 Option Groups allow profiles to subset the System Interfaces volume of IEEE Std 1003.1-200x by
 11614 collecting sets of related functions.

11615 POSIX_C_LANG_JUMP: Jump Functions

11616 *longjmp()*, *setjmp()*

11617 POSIX_C_LANG_MATH: Maths Library

11618 *acos()*, *acosf()*, *acosh()*, *acoshf()*, *acoshl()*, *acosl()*, *asin()*, *asinf()*, *asinh()*, *asinhf()*, *asinhl()*,
 11619 *asinl()*, *atan()*, *atan2()*, *atan2f()*, *atan2l()*, *atanf()*, *atanh()*, *atanhf()*, *atanhl()*, *atanl()*, *cabs()*,
 11620 *cabsf()*, *cabsl()*, *cacos()*, *cacosf()*, *cacosh()*, *cacoshf()*, *cacoshl()*, *cacosl()*, *carg()*, *cargf()*, *cargl()*,
 11621 *casin()*, *casinf()*, *casinh()*, *casinhf()*, *casinhl()*, *casinl()*, *catan()*, *catanf()*, *catanh()*, *catanhf()*,
 11622 *catanhl()*, *catanl()*, *cbrt()*, *cbrtf()*, *cbrtl()*, *ccos()*, *ccosf()*, *ccosh()*, *ccoshf()*, *ccoshl()*, *ccosl()*,
 11623 *ceil()*, *ceilf()*, *ceill()*, *cexp()*, *cexpf()*, *cexpl()*, *cimag()*, *cimagf()*, *cimagl()*, *clog()*, *clogf()*, *clogl()*,
 11624 *conj()*, *conjf()*, *conjl()*, *copysign()*, *copysignf()*, *copysignl()*, *cos()*, *cosf()*, *cosh()*, *coshf()*,
 11625 *coshl()*, *cosl()*, *cpow()*, *cpowf()*, *cpowl()*, *cproj()*, *cprojf()*, *cprojl()*, *creal()*, *crealf()*, *creall()*,
 11626 *csin()*, *csinf()*, *csinh()*, *csinhf()*, *csinhl()*, *csinl()*, *csqrt()*, *csqrtf()*, *csqrtl()*, *ctan()*, *ctanf()*,
 11627 *ctanh()*, *ctanhf()*, *ctanhl()*, *ctanl()*, *erf()*, *erfc()*, *erfcf()*, *erfcl()*, *erff()*, *erfl()*, *exp()*, *exp2()*,
 11628 *exp2f()*, *exp2l()*, *expf()*, *expl()*, *expm1()*, *expm1f()*, *expm1l()*, *fabs()*, *fabsf()*, *fabsl()*, *fdim()*,
 11629 *fdimf()*, *fdiml()*, *floor()*, *floorf()*, *floorl()*, *fma()*, *fmaf()*, *fmal()*, *fmax()*, *fmaxf()*, *fmaxl()*, *fmin()*,
 11630 *fminf()*, *fminl()*, *fmod()*, *fmodf()*, *fmodl()*, *fpclassify()*, *frexp()*, *frexpf()*, *frexpl()*, *hypot()*,
 11631 *hypotf()*, *hypotl()*, *ilogb()*, *ilogbf()*, *ilogbl()*, *isfinite()*, *isgreater()*, *isgreaterequal()*, *isinf()*,
 11632 *isless()*, *islessequal()*, *islessgreater()*, *isnan()*, *isnormal()*, *isunordered()*, *ldexp()*, *ldexpf()*,
 11633 *ldexpl()*, *lgamma()*, *lgammaf()*, *lgammal()*, *llrint()*, *llrintf()*, *llrintl()*, *llround()*, *llroundf()*,
 11634 *llroundl()*, *log()*, *log10()*, *log10f()*, *log10l()*, *log1p()*, *log1pf()*, *log1pl()*, *log2()*, *log2f()*, *log2l()*,
 11635 *logb()*, *logbf()*, *logbl()*, *logf()*, *logl()*, *lrint()*, *lrintf()*, *lrintl()*, *lround()*, *lroundf()*, *lroundl()*,
 11636 *modf()*, *modff()*, *modfl()*, *nan()*, *nanf()*, *nanl()*, *nearbyint()*, *nearbyintf()*, *nearbyintl()*,
 11637 *nextafter()*, *nextafterf()*, *nextafterl()*, *nexttoward()*, *nexttowardf()*, *nexttowardl()*, *pow()*, *powf()*,
 11638 *powl()*, *remainder()*, *remainderf()*, *remainderl()*, *remquo()*, *remquof()*, *remquol()*, *rint()*, *rintf()*,
 11639 *rintl()*, *round()*, *roundf()*, *roundl()*, *scalbln()*, *scalblnf()*, *scalblnl()*, *scalbn()*, *scalbnf()*, *scalbnl()*,
 11640 *signbit()*, *sin()*, *sinf()*, *sinh()*, *sinhf()*, *sinhl()*, *sinl()*, *sqrt()*, *sqrtf()*, *sqrtl()*, *tan()*, *tanf()*,

11641

11642 4. These are equivalent to the Units of Functionality from IEEE Std 1003.13-1998.

11643	<i>tanh()</i> , <i>tanhf()</i> , <i>tanhL()</i> , <i>tanl()</i> , <i>tgamma()</i> , <i>tgammaf()</i> , <i>tgammaL()</i> , <i>trunc()</i> , <i>truncf()</i> , <i>truncL()</i>
11644	POSIX_C_LANG_SUPPORT: General ISO C Library
11645	<i>abs()</i> , <i>asctime()</i> , <i>atof()</i> , <i>atoi()</i> , <i>atol()</i> , <i>bsearch()</i> , <i>calloc()</i> , <i>ctime()</i> , <i>difftime()</i> , <i>div()</i> ,
11646	<i>feclearexcept()</i> , <i>fegetenv()</i> , <i>fegetexceptflag()</i> , <i>fegetround()</i> , <i>fehldexcept()</i> , <i>feraiseexcept()</i> ,
11647	<i>fesetenv()</i> , <i>fesetexceptflag()</i> , <i>fesetround()</i> , <i>fetestexcept()</i> , <i>feupdateenv()</i> , <i>free()</i> , <i>gmtime()</i> ,
11648	<i>imaxabs()</i> , <i>imaxdiv()</i> , <i>isalnum()</i> , <i>isalpha()</i> , <i>isblank()</i> , <i>iscntrl()</i> , <i>isdigit()</i> , <i>isgraph()</i> , <i>islower()</i> ,
11649	<i>isprint()</i> , <i>ispunct()</i> , <i>isspace()</i> , <i>isupper()</i> , <i>isxdigit()</i> , <i>labs()</i> , <i>ldiv()</i> , <i>llabs()</i> , <i>lldiv()</i> , <i>localeconv()</i> ,
11650	<i>localtime()</i> , <i>malloc()</i> , <i>memchr()</i> , <i>memcmp()</i> , <i>memcpy()</i> , <i>memmove()</i> , <i>memset()</i> , <i>mktime()</i> ,
11651	<i>qsort()</i> , <i>rand()</i> , <i>realloc()</i> , <i>setlocale()</i> , <i>snprintf()</i> , <i>sprintf()</i> , <i>srand()</i> , <i>sscanf()</i> , <i>strcat()</i> , <i>strchr()</i> ,
11652	<i>strcmp()</i> , <i>strcoll()</i> , <i>strcpy()</i> , <i>strcspn()</i> , <i>strerror()</i> , <i>strftime()</i> , <i>strlen()</i> , <i>strncat()</i> , <i>strncpy()</i> ,
11653	<i>strncpy()</i> , <i>strpbrk()</i> , <i>strrchr()</i> , <i>strspn()</i> , <i>strstr()</i> , <i>strtod()</i> , <i>strtof()</i> , <i>strtoimax()</i> , <i>strtok()</i> , <i>strtol()</i> ,
11654	<i>strtold()</i> , <i>strtoll()</i> , <i>strtol()</i> , <i>strtol()</i> , <i>strtol()</i> , <i>strtol()</i> , <i>strtol()</i> , <i>strtol()</i> , <i>strtol()</i> , <i>strtol()</i> , <i>strtol()</i> ,
11655	<i>tzname</i> , <i>tzset()</i> , <i>va_arg()</i> , <i>va_copy()</i> , <i>va_end()</i> , <i>va_start()</i> , <i>vsprintf()</i> , <i>vsprintf()</i> , <i>vsscanf()</i>
11656	POSIX_C_LANG_SUPPORT_R: Thread-Safe General ISO C Library
11657	<i>asctime_r()</i> , <i>ctime_r()</i> , <i>gmtime_r()</i> , <i>localtime_r()</i> , <i>rand_r()</i> , <i>strerror_r()</i> , <i>strtok_r()</i>
11658	POSIX_C_LANG_WIDE_CHAR: Wide-Character ISO C Library
11659	<i>btowc()</i> , <i>iswalnum()</i> , <i>iswalnum()</i> , <i>iswalpha()</i> , <i>iswblank()</i> , <i>iswcntrl()</i> , <i>iswctype()</i> , <i>iswdigit()</i> , <i>iswgraph()</i> ,
11660	<i>iswlower()</i> , <i>iswprint()</i> , <i>iswpunct()</i> , <i>iswspace()</i> , <i>iswupper()</i> , <i>iswxdigit()</i> , <i>mblen()</i> , <i>mbrlen()</i> ,
11661	<i>mbrtowc()</i> , <i>mbsinit()</i> , <i>mbsrtowcs()</i> , <i>mbstowcs()</i> , <i>mbtowc()</i> , <i>swprintf()</i> , <i>swscanf()</i> , <i>towctrans()</i> ,
11662	<i>towlower()</i> , <i>towupper()</i> , <i>vswprintf()</i> , <i>vswscanf()</i> , <i>wcrtomb()</i> , <i>wcscat()</i> , <i>wcschr()</i> , <i>wcscmp()</i> ,
11663	<i>wcscoll()</i> , <i>wcscpy()</i> , <i>wcscspn()</i> , <i>wcsftime()</i> , <i>wcslen()</i> , <i>wcsncat()</i> , <i>wcsncmp()</i> , <i>wcsncpy()</i> ,
11664	<i>wcspbrk()</i> , <i>wcsrchr()</i> , <i>wcstombs()</i> , <i>wcsspn()</i> , <i>wcsstr()</i> , <i>wcstod()</i> , <i>wcstof()</i> , <i>wcstoimax()</i> ,
11665	<i>wcstok()</i> , <i>wcstol()</i> , <i>wcstold()</i> , <i>wcstoll()</i> , <i>wcstombs()</i> , <i>wcstoul()</i> , <i>wcstoull()</i> , <i>wcstoumax()</i> ,
11666	<i>wcsxfrm()</i> , <i>wctob()</i> , <i>wctomb()</i> , <i>wctrans()</i> , <i>wctype()</i> , <i>wmemchr()</i> , <i>wmemcmp()</i> , <i>wmemcpy()</i> ,
11667	<i>wmemmove()</i> , <i>wmemset()</i>
11668	POSIX_C_LIB_EXT: General C Library Extension
11669	<i>fnmatch()</i> , <i>getopt()</i> , <i>optarg</i> , <i>opterr</i> , <i>optind</i> , <i>optopt</i>
11670	POSIX_DEVICE_IO: Device Input and Output
11671	<i>FD_CLR()</i> , <i>FD_ISSET()</i> , <i>FD_SET()</i> , <i>FD_ZERO()</i> , <i>clearerr()</i> , <i>close()</i> , <i>fclose()</i> , <i>fdopen()</i> , <i>feof()</i> ,
11672	<i>ferror()</i> , <i>fflush()</i> , <i>fgetc()</i> , <i>fgets()</i> , <i>fileno()</i> , <i>fopen()</i> , <i>fprintf()</i> , <i>fputc()</i> , <i>fputs()</i> , <i>fread()</i> , <i>freopen()</i> ,
11673	<i>fscanf()</i> , <i>fwrite()</i> , <i>getc()</i> , <i>getchar()</i> , <i>gets()</i> , <i>open()</i> , <i>perror()</i> , <i>printf()</i> , <i>pselect()</i> , <i>putc()</i> , <i>putchar()</i> ,
11674	<i>puts()</i> , <i>read()</i> , <i>scanf()</i> , <i>select()</i> , <i>setbuf()</i> , <i>setvbuf()</i> , <i>stderr</i> , <i>stdin</i> , <i>stdout</i> , <i>ungetc()</i> , <i>vfprintf()</i> ,
11675	<i>vscanf()</i> , <i>vprintf()</i> , <i>vscanf()</i> , <i>write()</i>
11676	POSIX_DEVICE_SPECIFIC: General Terminal
11677	<i>cfgetispeed()</i> , <i>cfgetospeed()</i> , <i>cfsetispeed()</i> , <i>cfsetospeed()</i> , <i>ctermid()</i> , <i>isatty()</i> , <i>tcdrain()</i> , <i>tclflow()</i> ,
11678	<i>tcflush()</i> , <i>tcgetattr()</i> , <i>tcsendbreak()</i> , <i>tcsetattr()</i> , <i>ttyname()</i>
11679	POSIX_DEVICE_SPECIFIC_R: Thread-Safe General Terminal
11680	<i>ttyname_r()</i>
11681	POSIX_FD_MGMT: File Descriptor Management
11682	<i>dup()</i> , <i>dup2()</i> , <i>fcntl()</i> , <i>fgetpos()</i> , <i>fseek()</i> , <i>fseeko()</i> , <i>fsetpos()</i> , <i>ftell()</i> , <i>ftello()</i> , <i>ftruncate()</i> , <i>lseek()</i> ,
11683	<i>rewind()</i>
11684	POSIX_FIFO: FIFO
11685	<i>mkfifo()</i>
11686	POSIX_FILE_ATTRIBUTES: File Attributes
11687	<i>chmod()</i> , <i>chown()</i> , <i>fchmod()</i> , <i>fchown()</i> , <i>umask()</i>
11688	POSIX_FILE_LOCKING: Thread-Safe Stdio Locking
11689	<i>flockfile()</i> , <i>ftrylockfile()</i> , <i>funlockfile()</i> , <i>getc_unlocked()</i> , <i>getchar_unlocked()</i> , <i>putc_unlocked()</i> ,

11690	<i>putchar_unlocked()</i>	
11691	POSIX_FILE_SYSTEM: File System	
11692	<i>access(), chdir(), closedir(), creat(), fpathconf(), fstat(), getcwd(), link(), mkdir(), opendir(),</i>	
11693	<i>pathconf(), readdir(), remove(), rename(), rewinddir(), rmdir(), stat(), tmpfile(), tmpnam(),</i>	
11694	<i>unlink(), utime()</i>	
11695	POSIX_FILE_SYSTEM_EXT: File System Extensions	
11696	<i>glob(), globfree()</i>	
11697	POSIX_FILE_SYSTEM_R: Thread-Safe File System	
11698	<i>readdir_r()</i>	
11699	POSIX_JOB_CONTROL: Job Control	
11700	<i>setpgid(), tcgetpgrp(), tcsetpgrp()</i>	
11701	POSIX_MULTI_PROCESS: Multiple Processes	
11702	<i>_Exit(), _exit(), assert(), atexit(), clock(), execl(), execl(), execlp(), execv(), execve(), execvp(),</i>	
11703	<i>exit(), fork(), getpgrp(), getpid(), getppid(), setsid(), sleep(), times(), wait(), waitpid()</i>	
11704	POSIX_NETWORKING: Networking	
11705	<i>accept(), bind(), connect(), endhostent(), endnetent(), endprotoent(), endservent(),</i>	
11706	<i>freeaddrinfo(), gai_strerror(), getaddrinfo(), gethostbyaddr(), gethostbyname(), gethostent(),</i>	
11707	<i>gethostname(), getnameinfo(), getnetbyaddr(), getnetbyname(), getnetent(), getpeername(),</i>	
11708	<i>getprotobyname(), getprotobynumber(), getprotoent(), getservbyname(), getservbyport(),</i>	
11709	<i>getservent(), getsockname(), getsockopt(), h_errno, htonl(), htons(), if_freenameindex(),</i>	
11710	<i>if_indextoname(), if_nameindex(), if_nametoindex(), inet_addr(), inet_ntoa(), inet_ntop(),</i>	
11711	<i>inet_pton(), listen(), ntohl(), ntohs(), recv(), recvfrom(), recvmsg(), send(), sendmsg(), sendto(),</i>	
11712	<i>sethostent(), setnetent(), setprotoent(), setservent(), setsockopt(), shutdown(), socket(),</i>	
11713	<i>socketpair()</i>	
11714	POSIX_PIPE: Pipe	
11715	<i>pipe()</i>	
11716	POSIX_REGEX: Regular Expressions	
11717	<i>regcomp(), regerror(), regexexec(), regfree()</i>	
11718	POSIX_SHELL_FUNC: Shell and Utilities	
11719	<i>pclose(), popen(), system(), wordexp(), wordfree()</i>	
11720	POSIX_SIGNALS: Signal	
11721	<i>abort(), alarm(), kill(), pause(), raise(), sigaction(), sigaddset(), sigdelset(), sigemptyset(),</i>	
11722	<i>sigfillset(), sigismember(), signal(), sigpending(), sigprocmask(), sigsuspend(), sigwait()</i>	
11723	POSIX_SIGNAL_JUMP: Signal Jump Functions	
11724	<i>siglongjmp(), sigsetjmp()</i>	
11725	POSIX_SINGLE_PROCESS: Single Process	
11726	<i>confstr(), environ, errno, getenv(), setenv(), sysconf(), uname(), unsetenv()</i>	
11727	POSIX_SYMBOLIC_LINKS: Symbolic Links	
11728	<i>lstat(), readlink(), symlink()</i>	
11729	POSIX_SYSTEM_DATABASE: System Database	
11730	<i>getgrgid(), getgrnam(), getpwnam(), getpwuid()</i>	
11731	POSIX_SYSTEM_DATABASE_R: Thread-Safe System Database	
11732	<i>getgrgid_r(), getgrnam_r(), getpwnam_r(), getpwuid_r()</i>	

11733	POSIX_USER_GROUPS: User and Group	
11734	<i>getegid(), geteuid(), getgid(), getgroups(), getlogin(), getuid(), setegid(), seteuid(), setgid(),</i>	
11735	<i>setuid()</i>	
11736	POSIX_USER_GROUPS_R: Thread-Safe User and Group	
11737	<i>getlogin_r()</i>	
11738	POSIX_WIDE_CHAR_DEVICE_IO: Device Input and Output	
11739	<i>fgetwc(), fgetws(), fputwc(), fputws(), fwide(), fwprintf(), fwscanf(), getwc(), getwchar(),</i>	
11740	<i>putwc(), putwchar(), ungetwc(), vfwprintf(), vfwscanf(), vwprintf(), vwscanf(), wprintf(),</i>	
11741	<i>wscanf()</i>	
11742	XSI_C_LANG_SUPPORT: XSI General C Library	
11743	<i>_tolower(), _toupper(), a64l(), daylight(), drand48(), erand48(), ffs(), getcontext(), getdate(),</i>	
11744	<i>getsubopt(), hcreate(), hdestroy(), hsearch(), iconv(), iconv_close(), iconv_open(), initstate(),</i>	
11745	<i>insque(), isascii(), jrand48(), l64a(), lcong48(), lfind(), lrand48(), lsearch(), makecontext(),</i>	
11746	<i>memccpy(), mrand48(), nrand48(), random(), remque(), seed48(), setcontext(), setstate(),</i>	
11747	<i>signgam(), srand48(), srandom(), strcasecmp(), strdup(), strfmon(), strncasecmp(), strptime(),</i>	
11748	<i>swab(), swapcontext(), tdelete(), tfind(), timezone(), toascii(), tsearch(), twalk()</i>	
11749	XSI_DBM: XSI Database Management	
11750	<i>dbm_clearerr(), dbm_close(), dbm_delete(), dbm_error(), dbm_fetch(), dbm_firstkey(),</i>	
11751	<i>dbm_nextkey(), dbm_open(), dbm_store()</i>	
11752	XSI_DEVICE_IO: XSI Device Input and Output	
11753	<i>fntmsg(), poll(), pread(), pwrite(), readv(), writev()</i>	
11754	XSI_DEVICE_SPECIFIC: XSI General Terminal	
11755	<i>grantpt(), posix_openpt(), ptsname(), unlockpt()</i>	
11756	XSI_DYNAMIC_LINKING: XSI Dynamic Linking	
11757	<i>dlclose(), dlerror(), dlopen(), dlsym()</i>	
11758	XSI_FD_MGMT: XSI File Descriptor Management	
11759	<i>truncate()</i>	
11760	XSI_FILE_SYSTEM: XSI File System	
11761	<i>basename(), dirname(), fchdir(), fstatvfs(), ftw(), getwd(), lchown(), lockf(), mknod(),</i>	
11762	<i>mkstemp(), mktemp(), nftw(), realpath(), seekdir(), statvfs(), sync(), telldir(), tmpnam(),</i>	
11763	<i>utimes()</i>	
11764	XSI_I18N: XSI Internationalization	
11765	<i>catclose(), catgets(), catopen(), nl_langinfo()</i>	
11766	XSI_IPC: XSI Interprocess Communication	
11767	<i>ftok(), msgctl(), msgget(), msgrcv(), msgsnd(), semctl(), semget(), semop(), shmat(), shmctl(),</i>	
11768	<i>shmdt(), shmget()</i>	
11769	XSI_JOB_CONTROL: XSI Job Control	
11770	<i>tcgetsid()</i>	
11771	XSI_JUMP: XSI Jump Functions	
11772	<i>_longjmp(), _setjmp()</i>	
11773	XSI_MATH: XSI Maths Library	
11774	<i>j0(), j1(), jn(), scalb(), y0(), y1(), yn()</i>	
11775	XSI_MULTI_PROCESS: XSI Multiple Process	
11776	<i>getpgid(), getpriority(), getrlimit(), getrusage(), getsid(), nice(), setpgrp(), setpriority(),</i>	
11777	<i>setrlimit(), ulimit(), usleep(), vfork(), waitid()</i>	

11778	XSI_SIGNALS: XSI Signal	
11779	<i>bsd_signal()</i> , <i>killpg()</i> , <i>sigaltstack()</i> , <i>sighold()</i> , <i>sigignore()</i> , <i>siginterrupt()</i> , <i>sigpause()</i> , <i>sigrelse()</i> ,	
11780	<i>sigset()</i> , <i>ualarm()</i>	
11781	XSI_SINGLE_PROCESS: XSI Single Process	
11782	<i>ftime()</i> , <i>gethostid()</i> , <i>gettimeofday()</i> , <i>putenv()</i>	
11783	XSI_SYSTEM_DATABASE: XSI System Database	
11784	<i>endpwent()</i> , <i>getpwent()</i> , <i>setpwent()</i>	
11785	XSI_SYSTEM_LOGGING: XSI System Logging	
11786	<i>closelog()</i> , <i>openlog()</i> , <i>setlogmask()</i> , <i>syslog()</i>	
11787	XSI_THREAD_MUTEX_EXT: XSI Thread Mutex Extensions	
11788	<i>pthread_mutexattr_gettype()</i> , <i>pthread_mutexattr_settype()</i>	
11789	XSI_THREADS_EXT: XSI Threads Extensions	
11790	<i>pthread_attr_getguardsize()</i> , <i>pthread_attr_getstack()</i> , <i>pthread_attr_setguardsize()</i> ,	
11791	<i>pthread_attr_setstack()</i> , <i>pthread_getconcurrency()</i> , <i>pthread_setconcurrency()</i>	
11792	XSI_TIMERS: XSI Timers	
11793	<i>getitimer()</i> , <i>setitimer()</i>	
11794	XSI_USER_GROUPS: XSI User and Group	
11795	<i>endgrent()</i> , <i>endutxent()</i> , <i>getgrent()</i> , <i>getutxent()</i> , <i>getutxid()</i> , <i>getutxline()</i> , <i>pututxline()</i> ,	
11796	<i>setgrent()</i> , <i>setregid()</i> , <i>setreuid()</i> , <i>setutxent()</i>	
11797	XSI_WIDE_CHAR: XSI Wide-Character Library	
11798	<i>wcswcs()</i> , <i>wcswidth()</i> , <i>wcwidth()</i>	

