

1 / *Rationale (Informative)*

2 **Part A:**

3 **Base Definitions**

4 *The Open Group*



# Rationale for Base Definitions

5

## 6 A.1 Introduction

### 7 A.1.1 Scope

8 IEEE Std. 1003.1-200x is one of a family of standards known as POSIX. The family of standards  
9 extends to many topics; IEEE Std. 1003.1-200x is known as POSIX.1 and consists of both  
10 operating system interfaces and shell and utilities.

#### 11 Scope of IEEE Std. 1003.1-200x

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

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

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

23 The majority of the obsolescent material from the existing POSIX.1 and POSIX.2 standards, and  
24 material marked LEGACY from The Open Group's Base specifications, has been removed in this  
25 revision.

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

- 27 • IEEE Std. 1003.1d-1999
- 28 • IEEE Std. 1003.1j-2000
- 29 • IEEE Std. 1003.1q-2000
- 30 • IEEE P1003.1a draft standard
- 31 • IEEE Std. 1003.2d-1994
- 32 • IEEE P1003.2b draft standard
- 33 • Selected parts of IEEE Std. 1003.1g-2000

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

- 39 • General terms related to sockets (clause 2.2.2)

- 40 • Socket concepts (clauses 5.1 through 5.3 inclusive)
- 41 • The *pselect()* function (clauses 6.2.2.1 and 6.2.3)
- 42 • The *isfdtype()* function (clause 5.4.8)
- 43 • The `<sys/select.h>` header (clause 6.2)

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

- 45 • Backward-compatibility

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

- 50 • Architecture and n-bit neutral

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

- 54 • Extensibility

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

#### 58 **POSIX.1 and the ISO C standard**

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

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

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

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

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

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

**75 A.1.2 Conformance**

76 See Section A.2 (on page 3317).

**77 A.1.3 Normative References**

78 There is no additional rationale for this section.

**79 A.1.4 Terminology**

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

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

**85 conformance document**

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

**89 implementation-defined**

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

**93 may**

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

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

**99 shall**

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

- 103 1. The *foo()* function shall return zero.
- 104 2. The *foo()* function returns zero.

105 are meant to be exactly equivalent.

**106 should**

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

**110 obsolescent**

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

118 **legacy**

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

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

132 **system documentation**

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

139 **undefined**

140 See *implementation-defined*.

141 **unspecified**

142 See *implementation-defined*.

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

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

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

156 The term *implementation-defined* implies requirements for documentation that are not  
157 required for *undefined* (or *unspecified*). Where there is no need for a conforming program to  
158 know the definition, the term *undefined* is used, even though *implementation-defined* could  
159 also have been used in this context. There could be a fourth term, specifying “this standard  
160 does not say what this does; it is acceptable to define it in an implementation, but it does not  
161 need to be documented”, and *undefined* would then be used very rarely for the few things  
162 for which any definition is not useful.

163 In many places IEEE Std. 1003.1-200x is silent about the behavior of some possible construct.  
164 For example, a variable may be defined for a specified range of values and behaviors are  
165 described for those values; nothing is said about what happens if the variable has any other

166 value. That kind of silence can imply an error in the standard, but it may also imply that the  
 167 standard was intentionally silent and that any behavior is permitted. There is a natural  
 168 tendency to infer that if the standard is silent, a behavior is prohibited. That is not the intent.  
 169 Silence is intended to be equivalent to the term *unspecified*.

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

## 172 A.1.5 Portability

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

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

### 181 A.1.5.1 Codes

182 The set of code includes codes for options defined in clause 2.1.6, Options, and the following  
 183 additional codes for other purposes:

184	CX	This margin code is used to denote extensions beyond the ISO C standard, and is used
185		in interfaces that are duplicated between IEEE Std. 1003.1-200x and the ISO C standard.
186	MAN	This margin code was used during the development of the drafts and should not be
187		present in the final published standard.
188	OB	This margin code is used to denote obsolescent behavior and thus flag a possible future
189		application portability warning.
190	OH	The Single UNIX Specification has historically tried to reduce the number of headers an
191		application has had to include when using a particular interface. Sometimes this was
192		fewer than the base standard, and hence a notation is used to flag which headers are
193		optional if you are using a system supporting the XSI extension.
194	PI	This is another code used in the XSI extension only. It is used to denote a possible
195		application portability warning related to behavior of an interface which may not be
196		consistent between all conformant systems.
197	UN	This is another code used in the XSI extension only. It is used to denote a possible
198		application portability warning related to possibly unsupported functionality.
199	XSI	This code is used to denote interfaces and facilities within interfaces only required on
200		systems supporting the XSI extension. This is introduced to support the Single UNIX
201		Specification.
202	XSR	This code is used to denote interfaces and facilities within interfaces only required on
203		systems supporting STREAMS. This is introduced to support the Single UNIX
204		Specification, although it is defined in a way so that it can standalone from the XSI
205		notation.

206 *A.1.5.2 Margin Code Notation*

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



## 212 A.2 Conformance

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

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

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

- 218 • The addition of profiling options, allowing both sub-profiling as per IEEE Std. 1003.13-1998,  
219 and larger profiles of options such as the XSI extension used by the Single UNIX  
220 Specification. In effect, it has profiled itself (that is, created a self-profile).
- 221 • The addition of a hierarchy of super-options for XSI; these were formerly known as *Feature*  
222 *Groups* in The Open Group System Interfaces and Headers, Issue 5 specification.
- 223 • Options from the ISO POSIX-2: 1993 standard are also now included as IEEE Std. 1003.1-200x  
224 merges the functionality from it.

### 225 A.2.1 Implementation Conformance

226 These definitions allow application developers to know what to depend on in an  
227 implementation.

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

#### 232 A.2.1.1 Requirements

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

#### 236 A.2.1.2 Documentation

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

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

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

251 Also, interrelated file activities, where multiple files and/or directories are updated, or where  
252 space is allocated or released in the file system structures, can leave inconsistencies in the  
253 relationship between data in the various files and directories, or in the file system itself. Such  
254 inconsistencies can break applications that expect updates to occur in a specific sequence, so that

255 updates in one place correspond with related updates in another place.

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

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

### 269 A.2.1.3 *POSIX Conformance*

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

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

276 To enable support for sub-profiling the following options were defined:

- 277 • `_POSIX_C_LANG_SUPPORT`
- 278 • `_POSIX_DEVICE_IO`
- 279 • `_POSIX_DEVICE_SPECIFIC`
- 280 • `_POSIX_FD_MGMT`
- 281 • `_POSIX_FIFO`
- 282 • `_POSIX_FILE_ATTRIBUTES`
- 283 • `_POSIX_FILE_SYSTEM`
- 284 • `_POSIX_MULTIPLE_PROCESS`
- 285 • `_POSIX_PIPE`
- 286 • `_POSIX_SIGNALS`
- 287 • `_POSIX_SINGLE_PROCESS`
- 288 • `_POSIX_SYSTEM_DATABASE`
- 289 • `_POSIX_USER_GROUPS`
- 290 • `_POSIX_NETWORKING`

291 These are all mandatory in the base standard.

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

- 295           • `_POSIX_JOB_CONTROL`
- 296           • `_POSIX_SAVED_IDS`
- 297           • `_POSIX_VDISABLE`

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

#### 302   A.2.1.4   *XSI Conformance*

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

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

#### 307   A.2.1.5   *Option Groups*

308           The concept of *Option Groups* is introduced to IEEE Std. 1003.1-200x to allow collections of  
309           related functions or options to be grouped together. This is used in two ways in  
310           IEEE Std. 1003.1-200x:

- 311           1. Firstly, for profiling, a set of *Profiling Option Groups* has been created to support subsetting  
312           of the system interfaces provided in IEEE Std. 1003.1-200x. The subsets used by  
313           IEEE Std. 1003.13-1998 were used as an initial model for those created.
- 314           2. Secondly, the *XSI Option Groups* have been created to allow super-options, collections of  
315           underlying options and related functions, to be collectively supported by XSI-conforming  
316           systems. These reflect the *Feature Groups* from The Open Group System Interfaces and  
317           Headers, Issue 5 specification.

#### 318   A.2.1.6   *Options*

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

### 322   A.2.2    **Application Conformance**

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

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

333 A.2.2.1 *Strictly Conforming POSIX Application*334 This definition is analogous to that of a ISO C standard *conforming program*.335 The major difference between a *Strictly Conforming POSIX Application* and a ISO C standard  
336 *strictly conforming program* is that the latter is not allowed to use features of POSIX that are not in  
337 the ISO C standard.338 A.2.2.2 *Conforming POSIX Application*

339 Examples of &lt;National Bodies&gt; include ANSI, BSI, and AFNOR.

340 A.2.2.3 *Conforming POSIX Application Using Extensions*341 Due to possible requirements for configuration or implementation characteristics in excess of the  
342 specifications in <limits.h> or related to the hardware (such as array size or file space), not every  
343 Conforming POSIX Application Using Extensions will run on every conforming  
344 implementation.345 A.2.2.4 *Strictly Conforming XSI Application*346 This is intended to be upwards-compatible with the definition of a Strictly Conforming POSIX  
347 Application, with the addition of the facilities and functionality included in the XSI extension.348 A.2.2.5 *Conforming XSI Application Using Extensions*349 Such applications may use extensions beyond the facilities defined by IEEE Std. 1003.1-200x  
350 including the XSI extension, but need to document the additional requirements.351 **A.2.3 Language-Dependent Services for the C Programming Language**352 POSIX.1 is, for historical reasons, both a specification of an operating system interface, shell and  
353 utilities, and a C binding for that specification. Efforts had been previously undertaken to  
354 generate a language-independent specification; however, that had failed, and the fact that the  
355 ISO C standard is the *de facto* primary language on POSIX and the UNIX system makes this a  
356 necessary and workable situation.

### 357 A.3 Definitions

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

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

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

371 Some terms in the following list of notes do not appear in POSIX.1; these are marked prefixed  
372 with an asterisk (\*). Many of them have been specifically excluded from POSIX.1 because they  
373 concern system administration, implementation, or other issues that are not specific to the  
374 programming interface. Those are marked with a reason, such as “implementation-defined”.

#### 375 **Appropriate Privileges**

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

#### 390 **Character**

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

395 **Clock Tick**

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

399 **Command**

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

412 The terms *tool* and *application* can be thought of as being synonymous with *utility* from the  
413 perspective of the operating system kernel. Tools, applications, and utilities historically have  
414 run, typically, in processes above the kernel level. Tools and utilities historically have been a part  
415 of the operating system non-kernel code and have performed system-related functions, such as  
416 listing directory contents, checking file systems, repairing file systems, or extracting system  
417 status information. Applications have not generally been a part of the operating system, and  
418 they perform non-system-related functions, such as word processing, architectural design,  
419 mechanical design, workstation publishing, or financial analysis. Utilities have most frequently  
420 been provided by the operating system distributor, applications by third-party software  
421 distributors, or by the users themselves. Nevertheless, IEEE Std. 1003.1-200x does not  
422 differentiate between tools, utilities, and applications when it comes to receiving services from  
423 the system, a shell, or the standard utilities. (For example, the *xargs* utility invokes another  
424 utility; it would be of fairly limited usefulness if the users could not run their own applications  
425 in place of the standard utilities.) Utilities are not applications in the sense that they are not  
426 themselves subject to the restrictions of IEEE Std. 1003.1-200x or any other standard—there is no  
427 requirement for *grep*, *stty*, or any of the utilities defined here to be any of the classes of  
428 conforming applications.

429 **Column Positions**

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

436 In IEEE Std. 1003.1-200x the term *column positions* has been defined to mean character—not  
437 byte—positions in input files (such as “column position 7 of the FORTRAN input”). Output files  
438 describe the column position in terms of the display width of the narrowest printable character  
439 in the character set, adjusted to fit the characteristics of the output device. It is very possible that  
440 *n* column positions will not be able to hold *n* characters in some character sets, unless all of those  
441 characters are of the narrowest width. It is assumed that the implementation is aware of the  
442 width of the various characters, deriving this information from the value of *LC\_CTYPE*, and thus

443 can determine how many column positions to allot for each character in those utilities where it is  
444 important.

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

#### 449 **Controlling Terminal**

450 The question of which of possibly several special files referring to the terminal is meant is not  
451 addressed in POSIX.1. The file name `/dev/tty` is a synonym for the controlling terminal  
452 associated with a process.

#### 453 **Device Number\***

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

#### 455 **Direct I/O**

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

#### 458 **Directory**

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

#### 463 **Directory Entry**

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

#### 466 **Display**

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

**482 Dot**

483 The symbolic name *dot* is carefully used in POSIX.1 to distinguish the working directory file  
484 name from a period or a decimal point.

**485 Dot-Dot**

486 Historical implementations permit the use of these file names without their special meanings.  
487 Such use precludes any meaningful use of these file names by a Conforming POSIX.1  
488 Application. Therefore, such use is considered an extension, the use of which makes an  
489 implementation non-conforming; see also Section A.4.9 (on page 3346).

**490 Epoch**

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

**495 FIFO Special File**

496 See *pipe* in **Pipe** (on page 3331).

**497 File**

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

**499 File Classes**

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

**503 File Name**

504 At the present time, the primary responsibility for truncating file names containing multi-byte  
505 characters must reside with the application. Some industry groups involved in  
506 internationalization believe that in the future the responsibility must reside with the kernel. For  
507 the moment, a clearer understanding of the implications of making the kernel responsible for  
508 truncation of multi-byte file names is needed.

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

**513 File System**

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



**519 Graphic Character**

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

**522 Group Database**

523 See **User Database** (on page 3340).

**524 Group File\***

525 Implementation-defined; see **User Database** (on page 3340).

**526 Historical Implementations\***

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

**529 Hosted Implementation\***

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

**537 Implementation\***

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

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

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

554 **Incomplete Path Name\***

555 Absolute path name has been adequately defined.

556 **Job Control**

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

559 While the job control facilities supplied by POSIX.1 can, in theory, support different types of  
560 interactive job control interfaces supplied by different types of shells, there is historically one  
561 particular interface that is most common (provided by BSD C Shell). This discussion describes  
562 that interface as a means of illustrating how the POSIX.1 job control facilities can be used.

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

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

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

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

582 *Implementing Job Control Shells*

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

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

589 The concept of *process groups* is used to provide this grouping. The shell places each job in a  
590 separate process group via the *setpgid()* function. To do this, the *setpgid()* function is invoked by  
591 the shell for each process in the job. It is actually useful to invoke *setpgid()* twice for each  
592 process: once in the child process, after calling *fork()* to create the process, but before calling one  
593 of the *exec* family of functions to begin execution of the program, and once in the parent shell  
594 process, after calling *fork()* to create the child. The redundant invocation avoids a race condition  
595 by ensuring that the child process is placed into the new process group before either the parent  
596 or the child relies on this being the case. The *process group ID* for the job is selected by the shell to  
597 be equal to the *process ID* of one of the processes in the job. Some shells choose to make one  
598 process in the job be the parent of the other processes in the job (if any). Other shells (for  
599 example, the C Shell) choose to make themselves the parent of all processes in the pipeline (job).

600 In order to support this latter case, the *setpgid()* function accepts a process group ID parameter  
601 since the correct process group ID cannot be inherited from the shell. The shell itself is  
602 considered to be a job and is the sole process in its own process group.

603 The shell also controls which job is currently in the foreground. A foreground and background  
604 job differ in two ways: the shell waits for a foreground command to complete (or stop) before  
605 continuing to read new commands, and the terminal I/O driver inhibits terminal access by  
606 background jobs (causing the processes to stop). Thus, the shell must work cooperatively with  
607 the terminal I/O driver and have a common understanding of which job is currently in the  
608 foreground. It is the user who decides which command should be currently in the foreground,  
609 and the user informs the shell via shell commands. The shell, in turn, informs the terminal I/O  
610 driver via the *tcsetpgrp()* function. This indicates to the terminal I/O driver the process group ID  
611 of the foreground process group (job). When the current foreground job either stops or  
612 terminates, the shell places itself in the foreground via *tcsetpgrp()* before prompting for  
613 additional commands. Note that when a job is created the new process group begins as a  
614 background process group. It requires an explicit act of the shell via *tcsetpgrp()* to move a  
615 process group (job) into the foreground.

616 When a process in a job stops or terminates, its parent (for example, the shell) receives  
617 synchronous notification by calling the *waitpid()* function with the WUNTRACED flag set.  
618 Asynchronous notification is also provided when the parent establishes a signal handler for  
619 SIGCHLD and does not specify the SA\_NOCLDSTOP flag. Usually all processes in a job stop as  
620 a unit since the terminal I/O driver always sends job control stop signals to all processes in the  
621 process group.

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

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

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

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

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

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

644 A subshell can be spawned to perform an interactive function (prompting the terminal for  
645 commands) or a non-interactive function (reading commands from a file). When operating non-  
646 interactively, the job control shell will refrain from performing the job control-specific actions  
647 described above. It will behave as a shell that does not support job control. For example, all *jobs*

648 will be left in the same process group as the shell, which itself remains in the process group  
649 established for it by its parent. This allows the shell and its children to be treated as a single job  
650 by a parent shell, and they can be affected as a unit by terminal keyboard signals.

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

### 663 *Implementing Job Control Applications*

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

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

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

### 689 *Implementing Job Control Systems*

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

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

696 POSIX.1 does not specify how controlling terminal access is affected by a user logging out (that  
697 is, by a controlling process terminating). 4.2 BSD uses the *vhangup()* function to prevent any  
698 access to the controlling terminal through file descriptors opened prior to logout. System V does  
699 not prevent controlling terminal access through file descriptors opened prior to logout (except  
700 for the case of the special file, */dev/tty*). Some implementations choose to make processes  
701 immune from job control after logout (that is, such processes are always treated as if in the  
702 foreground); other implementations continue to enforce foreground/background checks after  
703 logout. Therefore, a Conforming POSIX.1 Application should not attempt to access the  
704 controlling terminal after logout since such access is unreliable. If an implementation chooses to  
705 deny access to a controlling terminal after its controlling process exits, POSIX.1 requires a certain  
706 type of behavior (see **Controlling Terminal** (on page 3323)).

#### 707 **Kernel\***

708 See *system call*.

#### 709 **Library Routine\***

710 See *system call*.

#### 711 **Logical Device\***

712 Implementation-defined.

#### 713 **Map**

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

#### 716 **Memory-Resident**

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

#### 730 **Memory Object\***

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

736 **Mount Point\***

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

739 **Mounted File System\***

740 See *file system*.

741 **name**

742 There are no explicit limits in IEEE Std. 1003.1-200x on the sizes of names, words (see the  
743 definition of word in the Base Definitions volume of IEEE Std. 1003.1-200x ), lines, or other  
744 objects. However, other implicit limits do apply: shell script lines produced by many of the  
745 standard utilities cannot exceed {LINE\_MAX} and the sum of exported variables comes under the  
746 {ARG\_MAX} limit. Historical shells dynamically allocate memory for names and words and  
747 parse incoming lines a byte at a time. Lines cannot have an arbitrary {LINE\_MAX} limit because  
748 of historical practice, such as makefiles, where *make* removes the <newline> characters  
749 associated with the commands for a target and presents the shell with one very long line. The  
750 text on INPUT FILES in the Shell and Utilities volume of IEEE Std. 1003.1-200x, Section 1.11,  
751 Utility Description Defaults does allow a shell to run out of memory, but it cannot have arbitrary  
752 programming limits.

753 **Native Implementation\***

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

758 **Nice Value**

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

766 **Open File Description**

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

770 The following alternate names were discussed:

- 771 • For *open file description*:  
772 *open instance, file access description, open file information, and file access information.*
- 773 • For *file descriptor*:  
774 *file handle, file number* (c.f., *fileno()*). Some historical implementations use the term *file table*  
775 *entry.*

**776 Orphaned Process Group**

777 Historical implementations have a concept of an orphaned process, which is a process whose  
778 parent process has exited. When job control is in use, it is necessary to prevent processes from  
779 being stopped in response to interactions with the terminal after they no longer are controlled by  
780 a job control-cognizant program. Because signals generated by the terminal are sent to a process  
781 group and not to individual processes, and because a signal may be provoked by a process that  
782 is not orphaned, but sent to another process that is orphaned, it is necessary to define an  
783 orphaned process group. The definition assumes that a process group will be manipulated as a  
784 group and that the job control-cognizant process controlling the group is outside of the group  
785 and is the parent of at least one process in the group (so that state changes may be reported via  
786 *waitpid()*). Therefore, a group is considered to be controlled as long as at least one process in the  
787 group has a parent that is outside of the process group, but within the session.

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

**791 Page**

792 The term *page* is defined to support the description of the behavior of memory mapping for  
793 shared memory and memory mapped files, and the description of the behavior of process  
794 memory locking. It is not intended to imply that shared memory/file mapping and memory  
795 locking are applicable only to "paged" architectures. For the purposes of IEEE Std. 1003.1-200x,  
796 whatever the granularity on which an architecture supports mapping or locking is considered to  
797 be a "page". If an architecture cannot support the memory mapping or locking functions  
798 specified by IEEE Std. 1003.1-200x on any granularity, then these options will not be  
799 implemented on the architecture.

**800 Passwd File\***

801 Implementation-defined; see **User Database** (on page 3340).

**802 Parent Directory**

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

**809 Pipe**

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

**812 Portable File Name Character Set**

813 The encoding of this character set is not specified—specifically, ASCII is not required. But the  
814 implementation must provide a unique character code for each of the printable graphics  
815 specified by POSIX.1; see also Section A.4.5 (on page 3342).

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

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

**830 Regular File**

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

**834 Root Directory**

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

**837 Root File System\***

838 Implementation-defined.

**839 Root of a File System\***

840 Implementation-defined; see *mount point*.

**841 Seconds Since the Epoch**

842 Coordinated Universal Time uses the concept of leap seconds; at the time POSIX.1 was  
843 published, 14 leap seconds had been added since January 1, 1970. These 14 seconds are ignored  
844 to provide an easy and compatible method of computing time differences.

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

850 It is sufficient to require that applications be allowed to treat this time as if it represented the  
851 number of seconds between the referenced time and the Epoch. It is the responsibility of the  
852 vendor of the system, and the administrator of the system, to ensure that this value represents  
853 the number of seconds between the referenced time and the Epoch as closely as necessary for the



854 application being run on that system.

855 It is important that the interpretation of time names and *seconds since the Epoch* values be  
856 consistent across conforming systems; that is, it is important that all conforming systems  
857 interpret “536 457 599 seconds since the Epoch” as 59 seconds, 59 minutes, 23 hours 31 December  
858 1986, regardless of the accuracy of the system’s idea of the current time. The expression is given  
859 to assure a consistent interpretation, not to attempt to specify the calendar. The relationship  
860 between *tm\_yday* and the day of week, day of month, and month is presumed to be specified  
861 elsewhere and is not given in POSIX.1.

862 Consistent interpretation of *seconds since the Epoch* can be critical to certain types of distributed  
863 applications that rely on such timestamps to synchronize events. The accrual of leap seconds in  
864 a time standard is not predictable. The number of leap seconds since the Epoch will likely  
865 increase. POSIX.1 is more concerned about the synchronization of time between applications of  
866 astronomically short duration. These concerns are expected to become more critical in the future.

867 Note that *tm\_yday* is zero-based, not one-based, so the day number in the example above is 364.  
868 Note also that the division is an integer division (discarding remainder) as in the C language.

869 Note also that the meaning of *gmtime()*, *localtime()*, and *mktime()* is specified in terms of this  
870 expression. However, the ISO C standard computes *tm\_yday* from *tm\_mday*, *tm\_mon*, and  
871 *tm\_year* in *mktime()*. Because it is stated as a (bidirectional) relationship, not a function, and  
872 because the conversion between month-day-year and day-of-year dates is presumed well known  
873 and is also a relationship, this is not a problem.

874 Implementations that implement **time\_t** as a 32-bit integer will overflow in 2 038. POSIX.1 does  
875 not Specify the data size for **time\_t**.

876 See also **Epoch** (on page 3324).

## 877 **Signal**

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

## 880 **Superuser\***

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

## 883 **Supplementary Group ID**

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

- 891 1. Reword the description of *setgid()* to permit it to change the supplementary group IDs to  
892 reflect the new effective group ID. A problem with this is that it adds more “may”s to the  
893 wording and does not address the portability problems of this optional behavior.
- 894 2. Mandate the inclusion of the effective group ID in the supplementary set (giving  
895 {NGROUPS\_MAX} a minimum value of 1). This is the behavior of 4.4 BSD. In that system,  
896 the effective group ID is the first element of the array of supplementary group IDs (there is  
897 no separate copy stored, and changes to the effective group ID are made only in the

898 supplementary group set). By convention, the initial value of the effective group ID is  
 899 duplicated elsewhere in the array so that the initial value is not lost when executing a set-  
 900 group-ID program.

901 3. Change the definition of supplementary group ID to exclude the effective group ID and  
 902 specify that the effective group ID does not change the set of supplementary group IDs.  
 903 This is the behavior of 4.2 BSD, 4.3 BSD, and System V, Release 4.

904 4. Change the definition of supplementary group ID to exclude the effective group ID, and  
 905 require that *getgroups()* return the union of the effective group ID and the supplementary  
 906 group IDs.

907 5. Change the definition of {NGROUPS\_MAX} to be one more than the number of  
 908 supplementary group IDs, so it continues to be the number of values returned by  
 909 *getgroups()* and existing applications continue to work. This alternative is effectively the  
 910 same as the second (and might actually have the same implementation).

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

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

## 924 **Symbolic Link**

925 Many implementations associate no attributes, including ownership with symbolic links.  
 926 Security experts encouraged consideration for defining these attributes as optional.  
 927 Consideration was given to changing *utime()* to allow modification of the times for a symbolic  
 928 link, or as an alternative adding an *lutime()* interface. Modifications to *chown()* were also  
 929 considered: allow changing symbolic link ownership or alternatively adding *lchown()*. As a  
 930 result of the problems encountered in defining attributes for symbolic links (and interfaces to  
 931 access/modify those attributes) and since implementations exist that do not associate these  
 932 attributes with symbolic links, only the file type bits in the *st\_mode* member and the *st\_size*  
 933 member of the **stat** structure are required to be applicable to symbolic links.

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

942 In the case of symbolic links, IEEE Std. 1003.1-200x states that a trailing slash is considered to be  
 943 the final component of a path name rather than the path name component that preceded it. This  
 944 is the behavior of historical implementations. For example, for **/a/b** and **/a/b/**, if **/a/b** is a symbolic

945 link to a directory, then `/a/b` refers to the symbolic link, and `/a/b/` is the same as `/a/b/.`, which is the  
946 directory to which the symbolic link points.

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

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

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

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

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

- 974 1. Symbolic links specified to system calls that take file name arguments
- 975 2. Symbolic links specified as command line file name arguments to utilities that are not  
976 performing a traversal of a file hierarchy
- 977 3. Symbolic links referencing files not of type directory, specified to utilities that are  
978 performing a traversal of a file hierarchy
- 979 4. Symbolic links referencing files of type directory, specified to utilities that are performing a  
980 traversal of a file hierarchy

#### 981 *First Domain*

982 The first domain is considered in earlier rationale.

#### 983 *Second Domain*

984 The reason this category is restricted to utilities that are not traversing the file hierarchy is that  
985 some standard utilities take an option that specifies a hierarchical traversal, but by default  
986 operate on the arguments themselves. Generally, users specifying the option for a file hierarchy  
987 traversal wish to operate on a single, physical hierarchy, and therefore symbolic links, which  
988 may reference files outside of the hierarchy, are ignored. For example, *chown owner file* is a  
989 different operation from the same command with the `-R` option specified. In this example, the  
990 behavior of the command *chown owner file* is described here, while the behavior of the command

- 991 *chown* **-R** *owner file* is described in the third and fourth domains.
- 992 The general rule is that the utilities in this category follow symbolic links named as arguments.
- 993 Exceptions in the second domain are:
- 994 • The *mv* and *rm* utilities do not follow symbolic links named as arguments, but respectively  
995 attempt to rename or delete them.
  - 996 • The *ls* utility is also an exception to this rule. For compatibility with historical systems, when  
997 the **-R** option is not specified, the *ls* utility follows symbolic links named as arguments if the  
998 **-L** option is specified or if the **-F**, **-d**, or **-l** options are not specified. (If the **-L** option is  
999 specified, *ls* always follows symbolic links; it is the only utility where the **-L** option affects its  
1000 behavior even though a tree walk is not being performed.)
- 1001 All other standard utilities, when not traversing a file hierarchy, always follow symbolic links  
1002 named as arguments.
- 1003 Historical practice is that the **-h** option is specified if standard utilities are to act upon symbolic  
1004 links instead of upon their targets. Examples of commands that have historically had a **-h** option  
1005 for this purpose are the *chgrp*, *chown*, *file*, and *test* utilities.
- 1006 *Third Domain*
- 1007 The third domain is symbolic links, referencing files not of type directory, specified to utilities  
1008 that are performing a traversal of a file hierarchy. (This includes symbolic links specified as  
1009 command line file name arguments or encountered during the traversal.)
- 1010 The intention of the Shell and Utilities volume of IEEE Std. 1003.1-200x is that the operation that  
1011 the utility is performing is applied to the symbolic link itself, if that operation is applicable to  
1012 symbolic links. The reason that the operation is not required is that symbolic links in some  
1013 systems do not have such attributes as a file owner, and therefore the *chown* operation would be  
1014 meaningless. If symbolic links on the system have an owner, it is the intention that the utility  
1015 *chown* cause the owner of the symbolic link to change. If symbolic links do not have an owner,  
1016 the symbolic link should be ignored. Specifically, by default, no change should be made to the  
1017 file referenced by the symbolic link.
- 1018 *Fourth Domain*
- 1019 The fourth domain is symbolic links referencing files of type directory, specified to utilities that  
1020 are performing a traversal of a file hierarchy. (This includes symbolic links specified as  
1021 command line file name arguments or encountered during the traversal.)
- 1022 All standard utilities do not, by default, indirect into the file hierarchy referenced by the  
1023 symbolic link. (The Shell and Utilities volume of IEEE Std. 1003.1-200x uses the informal term  
1024 *physical walk* to describe this case. The case where the utility does indirect through the symbolic  
1025 link is termed a *logical walk*.)
- 1026 There are three reasons for the default to a physical walk:
- 1027 1. With very few exceptions, a physical walk has been the historical default on UNIX systems  
1028 supporting symbolic links. Because some utilities (that is, *rm*) must default to a physical  
1029 walk, regardless, changing historical practice in this regard would be confusing to users  
1030 and needlessly incompatible.
  - 1031 2. For systems where symbolic links have the historical file attributes (that is, *owner*, *group*,  
1032 *mode*), defaulting to a logical traversal would require the addition of a new option to the  
1033 commands to modify the attributes of the link itself. This is painful and more complex  
1034 than the alternatives.

1035 3. There is a security issue with defaulting to a logical walk. Historically, the command  
1036 *chown -R user file* has been safe for the superuser because *setuid* and *setgid* bits were lost  
1037 when the ownership of the file was changed. If the walk were logical, changing ownership  
1038 would no longer be safe because a user might have inserted a symbolic link pointing to any  
1039 file in the tree. Again, this would necessitate the addition of an option to the commands  
1040 doing hierarchy traversal to not indirect through the symbolic links, and historical scripts  
1041 doing recursive walks would instantly become security problems. While this is mostly an  
1042 issue for system administrators, it is preferable to not have different defaults for different  
1043 classes of users.

1044 As consistently as possible, users may cause standard utilities performing a file hierarchy  
1045 traversal to follow any symbolic links named on the command line, regardless of the type of file  
1046 they reference, by specifying the **-H** (for half logical) option. This option is intended to make the  
1047 command line name space look like the logical name space.

1048 As consistently as possible, users may cause standard utilities performing a file hierarchy  
1049 traversal to follow any symbolic links named on the command line as well as any symbolic links  
1050 encountered during the traversal, regardless of the type of file they reference, by specifying the  
1051 **-L** (for logical) option. This option is intended to make the entire name space look like the  
1052 logical name space.

1053 For consistency, implementors are encouraged to use the **-P** (for physical) flag to specify the  
1054 physical walk in utilities that do logical walks by default for whatever reason. The only standard  
1055 utilities that require the **-P** option are *cd* and *pwd*; see the note below.

1056 When one or more of the **-H**, **-L**, and **-P** flags can be specified, the last one specified determines  
1057 the behavior of the utility. This permits users to alias commands so that the default behavior is a  
1058 logical walk and then override that behavior on the command line.

#### 1059 *Exceptions in the Third and Fourth Domains*

1060 The *ls* and *rm* utilities are exceptions to these rules. The *rm* utility never follows symbolic links  
1061 and does not support the **-H**, **-L**, or **-P** options. Some historical versions of *ls* always followed  
1062 symbolic links given on the command line whether the **-L** option was specified or not. Historical  
1063 versions of *ls* did not support the **-H** option. In IEEE Std. 1003.1-200x, the *ls* utility never follows  
1064 symbolic links unless one of the **-H** or **-L** options is specified. The *ls* utility does not support the  
1065 **-P** option.

1066 The Shell and Utilities volume of IEEE Std. 1003.1-200x requires that the standard utilities *ls*, *find*,  
1067 and *pax* detect infinite loops when doing logical walks; that is, a directory, or more commonly a  
1068 symbolic link, that refers to an ancestor in the current file hierarchy. If the file system itself is  
1069 corrupted, causing the infinite loop, it may be impossible to recover. Because *find* and *ls* are often  
1070 used in system administration and security applications, they should attempt to recover and  
1071 continue as best as they can. The *pax* utility should terminate because the archive it was creating  
1072 is by definition corrupted. Other, less vital, utilities should probably simply terminate as well.  
1073 Implementations are strongly encouraged to detect infinite loops in all utilities.

1074 Historical practice is shown in Table A-1 (on page 3338). The heading **SVID3** stands for the  
1075 Third Edition of the System V Interface Definition.

1076 Historically, several shells have had built-in versions of the *pwd* utility. In some of these shells,  
1077 *pwd* reported the physical path, and in others, the logical path. Implementations of the shell  
1078 corresponding to IEEE Std. 1003.1-200x must report the logical path by default. Earlier versions  
1079 of IEEE Std. 1003.1-200x did not require the *pwd* utility to be a built-in utility. Now that *pwd* is  
1080 required to set an environment variable in the current shell execution environment, it must be a  
1081 built-in utility.

1082 The *cd* command is required, by default, to treat the file name dot-dot logically. Implementors  
 1083 are required to support the **-P** flag in *cd* so that users can have their current environment  
 1084 handled physically. In 4.3 BSD, *chgrp* during tree traversal changed the group of the symbolic  
 1085 link, not the target. Symbolic links in 4.4 BSD do not have *owner*, *group*, *mode*, or other standard  
 1086 UNIX system file attributes.

1087 **Table A-1** Historical Practice for Symbolic Links

Utility	SVID3	4.3 BSD	4.4 BSD	POSIX	Comments
1088 <i>cd</i>				-L	Treat ". ." logically.
1089 <i>cd</i>				-P	". ." physically.
1090 <i>chgrp</i>			-H	-H	Follow command line symlinks.
1091 <i>chgrp</i>			-h	-L	Follow symlinks.
1092 <i>chgrp</i>	-h			-h	Affect the symlink.
1093 <i>chmod</i>				-h	Affect the symlink.
1094 <i>chmod</i>			-H	-H	Follow command line symlinks.
1095 <i>chmod</i>			-h	-L	Follow symlinks.
1096 <i>chown</i>			-H	-H	Follow command line symlinks.
1097 <i>chown</i>			-h	-L	Follow symlinks.
1098 <i>chown</i>	-h			-h	Affect the symlink.
1099 <i>cp</i>			-H	-H	Follow command line symlinks.
1100 <i>cp</i>			-h	-L	Follow symlinks.
1101 <i>cpio</i>	-L		-L		Follow symlinks.
1102 <i>du</i>			-H	-H	Follow command line symlinks.
1103 <i>du</i>			-h	-L	Follow symlinks.
1104 <i>file</i>	-h			-h	Affect the symlink.
1105 <i>find</i>			-H	-H	Follow command line symlinks.
1106 <i>find</i>			-h	-L	Follow symlinks.
1107 <i>find</i>	-follow		-follow		Follow symlinks.
1108 <i>ln</i>	-s	-s	-s	-s	Create a symbolic link.
1109 <i>ls</i>	-L	-L	-L	-L	Follow symlinks.
1110 <i>ls</i>				-H	Follow command line symlinks.
1111 <i>mv</i>					Operates on the symlink.
1112 <i>pax</i>			-H	-H	Follow command line symlinks.
1113 <i>pax</i>			-h	-L	Follow symlinks.
1114 <i>pwd</i>				-L	Printed path may contain symlinks.
1115 <i>pwd</i>				-P	Printed path will not contain symlinks.
1116 <i>rm</i>					Operates on the symlink.
1117 <i>tar</i>			-H		Follow command line symlinks.
1118 <i>tar</i>		-h	-h		Follow symlinks.
1119 <i>test</i>	-h		-h	-h	Affect the symlink.

### 1121 Synchronously-Generated Signal

1122 Those signals that may be generated synchronously include SIGABRT, SIGBUS, SIGILL, SIGFPE,  
 1123 SIGPIPE, and SIGSEGV.

**1124 System Call\***

1125 The distinction between a *system call* and a *library routine* is an implementation detail that may  
1126 differ between implementations and has thus been excluded from POSIX.1.

1127 See “Interface, Not Implementation” in the Introduction.

**1128 System Reboot**

1129 A *system reboot* is an event initiated by an unspecified circumstance that causes all processes  
1130 (other than special system processes) to be terminated in an implementation-defined manner,  
1131 after which any changes to the state and contents of files created or written to by a Conforming  
1132 POSIX.1 Application prior to the event are implementation-defined.

**1133 Synchronized I/O Data (and File) Integrity Completion**

1134 These terms specify that for synchronized read operations, pending writes must be successfully  
1135 completed before the read operation can complete. This is motivated by two circumstances.  
1136 Firstly, when synchronizing processes can access the same file, but not share common buffers  
1137 (such as for a remote file system), this requirement permits the reading process to guarantee that  
1138 it can read data written remotely. Secondly, having data written synchronously is insufficient to  
1139 guarantee the order with respect to a subsequent write by a reading process, and thus this extra  
1140 read semantic is necessary.

**1141 Text File**

1142 The term *text file* does not prevent the inclusion of control or other non-printable characters  
1143 (other than NUL). Therefore, standard utilities that list text files as inputs or outputs are either  
1144 able to process the special characters or they explicitly describe their limitations within their  
1145 individual descriptions. The definition of *text file* has caused controversy. The only difference  
1146 between text and binary files is that text files have lines of less than {LINE\_MAX} bytes, with no  
1147 NUL characters, each terminated by a <newline> character. The definition allows a file with a  
1148 single <newline> character, but not a totally empty file, to be called a text file. If a file ends with  
1149 an incomplete line it is not strictly a text file by this definition. The <newline> character referred  
1150 to in IEEE Std. 1003.1-200x is not some generic line separator, but a single character; files created  
1151 on systems where they use multiple characters for ends of lines are not portable to all  
1152 conforming systems without some translation process unspecified by IEEE Std. 1003.1-200x.

**1153 Thread**

1154 IEEE Std. 1003.1-200x defines a thread to be a flow of control within a process. Each thread has a  
1155 minimal amount of private state; most of the state associated with a process is shared among all  
1156 of the threads in the process. While most multi-thread extensions to POSIX have taken this  
1157 approach, others have made different decisions.

1158 **Note:** The choice to put threads within a process does not constrain implementations to  
1159 implement threads in that manner. However, all functions have to behave as though  
1160 threads share the indicated state information with the process from which they were  
1161 created.

1162 Threads need to share resources in order to cooperate. Memory has to be widely shared between  
1163 threads in order for the threads to cooperate at a fine level of granularity. Threads keep data  
1164 structures and the locks protecting those data structures in shared memory. For a data structure  
1165 to be usefully shared between threads, such structures should not refer to any data that can only  
1166 be interpreted meaningfully by a single thread. Thus, any system resources that might be  
1167 referred to in data structures need to be shared between all threads. File descriptors, path names,

1168 and pointers to stack variables are all things that programmers want to share between their  
1169 threads. Thus, the file descriptor table, the root directory, the current working directory, and the  
1170 address space have to be shared.

1171 Library implementations are possible as long as the effective behavior is as if system services  
1172 invoked by one thread do not suspend other threads. This may be difficult for some library  
1173 implementations on systems that do not provide asynchronous facilities.

1174 See Section B.2.9 (on page 3447) for additional rationale.

#### 1175 **Thread ID**

1176 See Section B.2.9.2 (on page 3463) for additional rationale.

#### 1177 **Thread-Safe Function**

1178 All functions required by IEEE Std. 1003.1-200x need to be thread-safe; see Section A.4.14 (on  
1179 page 3347) and Section B.2.9.1 (on page 3460) for additional rationale.

#### 1180 **User Database**

1181 There are no references in IEEE Std. 1003.1-200x to a *passwd file* or a *group file*, and there is no  
1182 requirement that the *group* or *passwd* databases be kept in files containing editable text. Many  
1183 large timesharing systems use *passwd* databases that are hashed for speed. Certain security  
1184 classifications prohibit certain information in the *passwd* database from being publicly readable.

1185 The term *encoded* is used instead of *encrypted* in order to avoid the implementation connotations  
1186 (such as reversibility or use of a particular algorithm) of the latter term.

1187 The *getgrent()*, *setgrent()*, *endgrent()*, *getpwent()*, *setpwent()*, and *endpwent()* functions are not  
1188 included as part of the base standard because they provide a linear database search capability  
1189 that is not generally useful (the *getpwuid()*, *getpwnam()*, *getgrgid()*, and *getgrnam()* functions are  
1190 provided for keyed lookup) and because in certain distributed systems, especially those with  
1191 different authentication domains, it may not be possible or desirable to provide an application  
1192 with the ability to browse the system databases indiscriminately. They are provided on XSI-  
1193 conformant systems due to their historical usage by many existing applications.

1194 A change from historical implementations is that the structures used by these functions have  
1195 fields of the types **gid\_t** and **uid\_t**, which are required to be defined in the **<sys/types.h>** header.  
1196 IEEE Std. 1003.1-200x requires implementations to ensure that these types are defined by  
1197 inclusion of **<grp.h>** and **<pwd.h>**, respectively, without imposing any name space pollution or  
1198 errors from redefinition of types.

1199 IEEE Std. 1003.1-200x is silent about the content of the strings containing user or group names.  
1200 These could be digit strings. IEEE Std. 1003.1-200x is also silent as to whether such digit strings  
1201 bear any relationship to the corresponding (numeric) user or group ID.

#### 1202 *Database Access*

1203 The thread-safe versions of the user and group database access functions return values in user-  
1204 supplied buffers instead of possibly using static data areas that may be overwritten by each call.



1205 **Virtual Processor\***

1206 The term *virtual processor* was chosen as a neutral term describing all kernel-level schedulable  
1207 entities, such as processes, Mach tasks, or lightweight processes. Implementing threads using  
1208 multiple processes as virtual processors, or implementing multiplexed threads above a virtual  
1209 processor layer, should be possible, provided some mechanism has also been implemented for  
1210 sharing state between processes or virtual processors. Many systems may also wish to provide  
1211 implementations of threads on systems providing “shared processes” or “variable-weight  
1212 processes”. It was felt that exposing such implementation details would severely limit the type  
1213 of systems upon which the threads interface could be supported and prevent certain types of  
1214 valid implementations. It was also determined that a virtual processor interface was out of the  
1215 scope of the Rationale (Informative) volume of IEEE Std. 1003.1-200x.

1216 **XSI**

1217 This is introduced to allow IEEE Std. 1003.1-200x to be adopted as an IEEE standard and an  
1218 Open Group Technical Standard, serving both the POSIX and the Single UNIX Specification in a  
1219 core set of volumes.

1220 The term *XSI* has been used for 10 years in connection with the XPG series and the first and  
1221 second versions of the base volumes of the Single UNIX Specification. The XSI margin code was  
1222 introduced to denote the extended or more restrictive semantics beyond POSIX that are  
1223 applicable to UNIX systems.

## 1224 **A.4 General Concepts**

### 1225 **A.4.1 Concurrent Execution**

1226 There is no additional rationale provided for this section.

### 1227 **A.4.2 Extended Security Controls**

1228 Allowing an implementation to define extended security controls enables the use of  
1229 IEEE Std. 1003.1-200x in environments that require different or more rigorous security than that  
1230 provided in POSIX.1. Extensions are allowed in two areas: privilege and file access permissions.  
1231 The semantics of these areas have been defined to permit extensions with reasonable, but not  
1232 exact, compatibility with all existing practices. For example, the elimination of the superuser  
1233 definition precludes identifying a process as privileged or not by virtue of its effective user ID.

### 1234 **A.4.3 File Access Permissions**

1235 A process should not try to anticipate the result of an attempt to access data by *a priori* use of  
1236 these rules. Rather, it should make the attempt to access data and examine the return value (and  
1237 possibly *errno* as well), or use *access()*. An implementation may include other security  
1238 mechanisms in addition to those specified in POSIX.1, and an access attempt may fail because of  
1239 those additional mechanisms, even though it would succeed according to the rules given in this  
1240 section. (For example, the user's security level might be lower than that of the object of the access  
1241 attempt.) The supplementary group IDs provide another reason for a process to not attempt to  
1242 anticipate the result of an access attempt.

### 1243 **A.4.4 File Hierarchy**

1244 Though the file hierarchy is commonly regarded to be a tree, POSIX.1 does not define it as such  
1245 for three reasons:

- 1246 1. Links may join branches.
- 1247 2. In some network implementations, there may be no single absolute root directory; see *path*  
1248 *name resolution*.
- 1249 3. With symbolic links, the file system need not be a tree or even a directed acyclic graph.

### 1250 **A.4.5 File Names**

1251 Historically, certain file names have been reserved. This list includes *core*, */etc/passwd*, and so  
1252 on. Portable applications should avoid these.

1253 Most historical implementations prohibit case folding in file names; that is, treating uppercase  
1254 and lowercase alphabetic characters as identical. However, some consider case folding desirable:

- 1255 • For user convenience
- 1256 • For ease-of-implementation of the POSIX.1 interface as a hosted system on some popular  
1257 operating systems

1258 Variants, such as maintaining case distinctions in file names, but ignoring them in comparisons,  
1259 have been suggested. Methods of allowing escaped characters of the case opposite the default  
1260 have been proposed.

1261 Many reasons have been expressed for not allowing case folding, including:

- 1262 • No solid evidence has been produced as to whether case-sensitivity or case-insensitivity is  
1263 more convenient for users.
- 1264 • Making case-insensitivity a POSIX.1 implementation option would be worse than either  
1265 having it or not having it, because:
- 1266 — More confusion would be caused among users.
- 1267 — Application developers would have to account for both cases in their code.
- 1268 — POSIX.1 implementors would still have other problems with native file systems, such as  
1269 short or otherwise constrained file names or path names, and the lack of hierarchical  
1270 directory structure.
- 1271 • Case folding is not easily defined in many European languages, both because many of them  
1272 use characters outside the USASCII alphabetic set, and because:
- 1273 — In Spanish, the digraph "ll" is considered to be a single letter, the capitalized form of  
1274 which may be either "Ll" or "LL", depending on context.
- 1275 — In French, the capitalized form of a letter with an accent may or may not retain the accent,  
1276 depending on the country in which it is written.
- 1277 — In German, the sharp ess may be represented as a single character resembling a Greek  
1278 beta (β) in lowercase, but as the digraph "SS" in uppercase.
- 1279 — In Greek, there are several lowercase forms of some letters; the one to use depends on its  
1280 position in the word. Arabic has similar rules.
- 1281 • Many East Asian languages, including Japanese, Chinese, and Korean, do not distinguish  
1282 case and are sometimes encoded in character sets that use more than one byte per character.
- 1283 • Multiple character codes may be used on the same machine simultaneously. There are  
1284 several ISO character sets for European alphabets. In Japan, several Japanese character codes  
1285 are commonly used together, sometimes even in file names; this is evidently also the case in  
1286 China. To handle case insensitivity, the kernel would have to at least be able to distinguish  
1287 for which character sets the concept made sense.
- 1288 • The file system implementation historically deals only with bytes, not with characters, except  
1289 for slash and the null byte.
- 1290 • The purpose of POSIX.1 is to standardize the common, existing definition, not to change it.  
1291 Mandating case-insensitivity would make all historical implementations non-standard.
- 1292 • Not only the interface, but also application programs would need to change, counter to the  
1293 purpose of having minimal changes to existing application code.
- 1294 • At least one of the original developers of the UNIX system has expressed objection in the  
1295 strongest terms to either requiring case-insensitivity or making it an option, mostly on the  
1296 basis that POSIX.1 should not hinder portability of application programs across related  
1297 implementations in order to allow compatibility with unrelated operating systems.
- 1298 Two proposals were entertained regarding case folding in file names:
- 1299 1. Remove all wording that previously permitted case folding.
- 1300 Rationale Case folding is inconsistent with portable file name character set definition  
1301 and file name definition (all characters except slash and null). No known  
1302 implementations allowing all characters except slash and null also do case  
1303 folding.

1304 2. Change “though this practice is not recommended:” to “although this practice is strongly  
1305 discouraged.”

1306 Rationale If case folding must be included in POSIX.1, the wording should be stronger  
1307 to discourage the practice.

1308 The consensus selected the first proposal. Otherwise, a portable application would have to  
1309 assume that case folding would occur when it was not wanted, but that it would not occur when  
1310 it was wanted.

#### 1311 **A.4.6 File Times Update**

1312 This section reflects the actions of historical implementations. The times are not updated  
1313 immediately, but are only marked for update by the functions. An implementation may update  
1314 these times immediately.

1315 The accuracy of the time update values is intentionally left unspecified so that systems can  
1316 control the bandwidth of a possible covert channel.

1317 The wording was carefully chosen to make it clear that there is no requirement that the  
1318 conformance document contain information that might incidentally affect file update times. Any  
1319 function that performs path name resolution might update several *st\_atime* fields. Functions  
1320 such as *getpwnam()* and *getgrnam()* might update the *st\_atime* field of some specific file or files. It  
1321 is intended that these are not required to be documented in the conformance document, but they  
1322 should appear in the system documentation.

#### 1323 **A.4.7 Measurement of Execution Time**

1324 The methods used to measure the execution time of processes and threads, and the precision of  
1325 these measurements, may vary considerably depending on the software architecture of the  
1326 implementation, and on the underlying hardware. Implementations can also make tradeoffs  
1327 between the scheduling overhead and the precision of the execution time measurements.  
1328 IEEE Std. 1003.1-200x does not impose any requirement on the accuracy of the execution time; it  
1329 instead specifies that the measurement mechanism and its precision are implementation-  
1330 defined.

#### 1331 **A.4.8 Memory Synchronization**

1332 In older multi-processors, access to memory by the processors was strictly multiplexed. This  
1333 meant that a processor executing program code interrogates or modifies memory in the order  
1334 specified by the code and that all the memory operation of all the processors in the system  
1335 appear to happen in some global order, though the operation histories of different processors are  
1336 interleaved arbitrarily. The memory operations of such machines are said to be sequentially  
1337 consistent. In this environment, threads can synchronize using ordinary memory operations. For  
1338 example, a producer thread and a consumer thread can synchronize access to a circular data  
1339 buffer as follows:

```

1340     int rdptr = 0;
1341     int wrptr = 0;
1342     data_t buf[BUFSIZE];

1343     Thread 1:
1344         while (work_to_do) {
1345             int next;

1346             buf[wrptr] = produce();
1347             next = (wrptr + 1) % BUFSIZE;
1348             while (rdptr == next)
1349                 ;
1350             wrptr = next;
1351         }

1352     Thread 2:
1353         while (work_to_do) {
1354             while (rdptr == wrptr)
1355                 ;
1356             consume(buf[rdptr]);
1357             rdptr = (rdptr + 1) % BUFSIZE;
1358         }

```

1359 In modern multi-processors, these conditions are relaxed to achieve greater performance. If one  
1360 processor stores values in location A and then location B, then other processors loading data  
1361 from location B and then location A may see the new value of B but the old value of A. The  
1362 memory operations of such machines are said to be weakly ordered. On these machines, the  
1363 circular buffer technique shown in the example will fail because the consumer may see the new  
1364 value of *wrptr* but the old value of the data in the buffer. In such machines, synchronization can  
1365 only be achieved through the use of special instructions that enforce an order on memory  
1366 operations. Most high-level language compilers only generate ordinary memory operations to  
1367 take advantage of the increased performance. They usually cannot determine when memory  
1368 operation order is important and generate the special ordering instructions. Instead, they rely on  
1369 the programmer to use synchronization primitives correctly to ensure that modifications to a  
1370 location in memory are ordered with respect to modifications and/or access to the same location  
1371 in other threads. Access to read-only data need not be synchronized. The resulting program is  
1372 said to be data race-free.

1373 Synchronization is still important even when accessing a single primitive variable (for example,  
1374 an integer). On machines where the integer may not be aligned to the bus data width or be larger  
1375 than the data width, a single memory load may require multiple memory cycles. This means  
1376 that it may be possible for some parts of the integer to have an old value while other parts have a  
1377 newer value. On some processor architectures this cannot happen, but portable programs cannot  
1378 rely on this.

1379 In summary, a portable multi-threaded program, or a multi-process program that shares  
1380 writable memory between processes, has to use the synchronization primitives to synchronize  
1381 data access. It cannot rely on modifications to memory being observed by other threads in the  
1382 order written in the program or even on modification of a single variable being seen atomically.

1383 Conforming applications may only use the functions listed to synchronize threads of control  
1384 with respect to memory access. There are many other candidates for functions that might also be  
1385 used. Examples are: signal sending and reception, or pipe writing and reading. In general, any  
1386 function that allows one thread of control to wait for an action caused by another thread of  
1387 control is a candidate. IEEE Std. 1003.1-200x does not require these additional functions to  
1388 synchronize memory access since this would imply the following:

- 1389 • All these functions would have to be recognized by advanced compilation systems so that
- 1390 memory operations and calls to these functions are not reordered by optimization.
- 1391 • All these functions would potentially have to have memory synchronization instructions
- 1392 added, depending on the particular machine.
- 1393 • The additional functions complicate the model of how memory is synchronized and make
- 1394 automatic data race detection techniques impractical.

1395 Formal definitions of the memory model were rejected as unreadable by the vast majority of  
 1396 programmers. In addition, most of the formal work in the literature has concentrated on the  
 1397 memory as provided by the hardware as opposed to the application programmer through the  
 1398 compiler and runtime system. It was believed that a simple statement intuitive to most  
 1399 programmers would be most effective. IEEE Std. 1003.1-200x defines functions that can be used  
 1400 to synchronize access to memory, but it leaves open exactly how one relates those functions to  
 1401 the semantics of each function as specified elsewhere in IEEE Std. 1003.1-200x.  
 1402 IEEE Std. 1003.1-200x also does not make a formal specification of the partial ordering in time  
 1403 that the functions can impose, as that is implied in the description of the semantics of each  
 1404 function. It simply states that the programmer has to ensure that modifications do not occur  
 1405 “simultaneously” with other access to a memory location.

#### 1406 A.4.9 Path Name Resolution

1407 It is necessary to differentiate between the definition of path name and the concept of path name  
 1408 resolution with respect to the handling of trailing slashes. By specifying the behavior here, it is  
 1409 not possible to provide an implementation that is conforming but extends all interfaces that  
 1410 handle path names to also handle strings that are not legal path names (because they have  
 1411 trailing slashes).

1412 Path names that end with one or more trailing slash characters must refer to directory paths.  
 1413 Previous versions of IEEE Std. 1003.1-200x were not specific about the distinction between  
 1414 trailing slashes on files and directories, and both were permitted.

1415 Two types of implementation have been prevalent; those that ignored trailing slash characters  
 1416 on all path names regardless, and those that only permitted them only on existing directories.

1417 IEEE Std. 1003.1-200x requires that a path name with a trailing slash character be treated as if it  
 1418 had a trailing " / . " everywhere.

1419 Note that this change does not break any portable applications; since there were two different  
 1420 types of implementation, no application could have portably depended on either behavior. This  
 1421 change does however require some implementations to be altered to remain compliant.  
 1422 Substantial discussion over a three-year period has shown that the benefits to application  
 1423 developers outweighs the disadvantages for some vendors.

1424 On a historical note, some early applications automatically appended a ' / ' to every path.  
 1425 Rather than fix the applications, the system implementation was modified to accept this  
 1426 behavior by ignoring any trailing slash.

1427 Each directory has exactly one parent directory which is represented by the name **dot-dot** in the  
 1428 first directory. No other directory, regardless of linkages established by symbolic links, is  
 1429 considered the parent directory by IEEE Std. 1003.1-200x.

1430 There are two general categories of interfaces involving path name resolution: those that follow  
 1431 the symbolic link, and those that do not. There are several exceptions to this rule; for example,  
 1432 *open(path, O\_CREAT | O\_EXCL)* will fail when *path* names a symbolic link. However, in all other  
 1433 situations, the *open()* function will follow the link.

1434 What the file name **dot-dot** refers to relative to the root directory is implementation-defined. In  
 1435 Version 7 it refers to the root directory itself; this is the behavior mentioned in  
 1436 IEEE Std. 1003.1-200x. In some networked systems the construction `././hostname/` is used to  
 1437 refer to the root directory of another host, and POSIX.1 permits this behavior.

1438 Other networked systems use the construct `//hostname` for the same purpose; that is, a double  
 1439 initial slash is used. There is a potential problem with existing applications that create full path  
 1440 names by taking a trunk and a relative path name and making them into a single string  
 1441 separated by `'/'`, because they can accidentally create networked path names when the trunk is  
 1442 `'/'`. This practice is not prohibited because such applications can be made to conform by  
 1443 simply changing to use `"/"` as a separator instead of `'/'`:

- 1444 • If the trunk is `'/'`, the full path name will begin with `"/"` (the initial `'/'` and the  
 1445 separator `"/"`). This is the same as `'/'`, which is what is desired. (This is the general case  
 1446 of making a relative path name into an absolute one by prefixing with `"/"` instead of `'/'`.)
- 1447 • If the trunk is `"/A"`, the result is `"/A/. . ."`; since non-leading sequences of two or more  
 1448 slashes are treated as a single slash, this is equivalent to the desired `"/A/. . ."`.
- 1449 • If the trunk is `"//A"`, the implementation-defined semantics will apply. (The multiple slash  
 1450 rule would apply.)

1451 Application developers should avoid generating path names that start with `"/"`.  
 1452 Implementations are strongly encouraged to avoid using this special interpretation since a  
 1453 number of applications currently do not follow this practice and may inadvertently generate  
 1454 `"/. . ."`.

1455 The term *root directory* is only defined in POSIX.1 relative to the process. In some  
 1456 implementations, there may be no absolute root directory. The initialization of the root directory  
 1457 of a process is implementation-defined.

#### 1458 **A.4.10 Process ID Reuse**

1459 There is no additional rationale provided for this section.

#### 1460 **A.4.11 Scheduling Policy**

1461 There is no additional rationale provided for this section.

#### 1462 **A.4.12 Seconds Since the Epoch**

1463 There is no additional rationale provided for this section.

#### 1464 **A.4.13 Semaphore**

1465 There is no additional rationale provided for this section.

#### 1466 **A.4.14 Thread-Safety**

1467 Where the interface of a function required by IEEE Std. 1003.1-200x precludes thread-safety, an  
 1468 alternate form that shall be thread-safe is provided. The names of these thread-safe forms are the  
 1469 same as the non-thread-safe forms with the addition of the suffix `"_r"`. The suffix `"_r"` is  
 1470 historical, where the `'r'` stood for "reentrant".

1471 In some cases, thread-safety is provided by restricting the arguments to an existing function.

1472 See also Section B.2.9.1 (on page 3460).

1473 **A.4.15 Utility**

1474           There is no additional rationale provided for this section.

1475 **A.4.16 Variable Assignment**

1476           There is no additional rationale provided for this section.



1477 **A.5 File Format Notation**

1478 The notation for spaces allows some flexibility for application output. Note that an empty  
 1479 character position in *format* represents one or more <blank> characters on the output (not *white*  
 1480 *space*, which can include <newline> characters). Therefore, another utility that reads that output  
 1481 as its input must be prepared to parse the data using *scanf()*, *awk*, and so on. The ' $\Delta$ ' character  
 1482 is used when exactly one <space> character is output.

1483 The treatment of integers and spaces is different from the *printf()* function in that they can be  
 1484 surrounded with <blank> characters. This was done so that, given a format such as:

```
1485 "%d\n", <foo>
```

1486 the implementation could use a *printf()* call such as:

```
1487 printf("%6d\n", foo);
```

1488 and still conform. This notation is thus somewhat like *scanf()* in addition to *printf()*.

1489 The *printf()* function was chosen as a model because most of the standard developers were  
 1490 familiar with it. One difference from the C function *printf()* is that the *l* and *h* conversion  
 1491 characters are not used. As expressed by the Shell and Utilities volume of IEEE Std. 1003.1-200x,  
 1492 there is no differentiation between decimal values for type **int**, type **long**, or type **short**. The  
 1493 specifications *%d* or *%i* should be interpreted as an arbitrary length sequence of digits. Also, no  
 1494 distinction is made between single precision and double precision numbers (**float** or **double** in  
 1495 C). These are simply referred to as floating point numbers.

1496 Many of the output descriptions in the Shell and Utilities volume of IEEE Std. 1003.1-200x use  
 1497 the term *line*, such as:

```
1498 "%s", <input line>
```

1499 Since the definition of *line* includes the trailing <newline> character already, there is no need to  
 1500 include a ' $\backslash n$ ' in the format; a double <newline> character would otherwise result.

**1501 A.6 Character Set****1502 A.6.1 Portable Character Set**

1503 The portable character set is listed in full so there is no dependency on the ISO/IEC 646:1991  
1504 standard (or historically ASCII) encoded character set, although the set is identical to the  
1505 characters defined in the International Reference version of the ISO/IEC 646:1991 standard.

1506 IEEE Std. 1003.1-200x poses no requirement that multiple character sets or codesets be  
1507 supported, leaving this as a marketing differentiation for implementors. Although multiple  
1508 charmap files are supported, it is the responsibility of the implementation to provide the file(s);  
1509 if only one is provided, only that one will be accessible using the *localedef -f* option.

1510 The statement about invariance in codesets for the portable character set is worded to avoid  
1511 precluding implementations where multiple incompatible codesets are available (for instance,  
1512 ASCII and EBCDIC). The standard utilities cannot be expected to produce predictable results if  
1513 they access portable characters that vary on the same implementation.

1514 Not all character sets need include the portable character set, but each locale must include it. For  
1515 example, a Japanese-based locale might be supported by a mixture of character sets: JIS X 0201  
1516 Roman (a Japanese version of the ISO/IEC 646:1991 standard), JIS X 0208, and JIS X 0201  
1517 Katakana. Not all of these character sets include the portable characters, but at least one does  
1518 (JIS X 0201 Roman).

**1519 A.6.2 Character Encoding**

1520 Encoding mechanisms based on single shifts, such as the EUC encoding used in some Asian and  
1521 other countries, can be supported via the current charmap mechanism. With single-shift  
1522 encoding, each character is preceded by a shift code (SS2 or SS3). A complete EUC code,  
1523 consisting of the portable character set (G0) and up to three additional character sets (G1, G2,  
1524 G3), can be described using the current charmap mechanism; the encoding for each character in  
1525 additional character sets G2 and G3 must then include their single-shift code. Other mechanisms  
1526 to support locales based on encoding mechanisms such as locking shift are not addressed by this  
1527 volume of IEEE Std. 1003.1-200x.

**1528 A.6.3 C Language Wide-Character Codes**

1529 There is no additional rationale for this section.

**1530 A.6.4 Character Set Description File****1531 A.6.4.1 State-Dependent Character Encodings**

1532 A requirement was considered that would force utilities to eliminate any redundant locking  
1533 shifts, but this was left as a quality of implementation issue.

1534 This change satisfies the following requirement from the ISO POSIX-2:1993 standard, Annex  
1535 H.1:

1536           The support of state-dependent (shift encoding) character sets should be addressed fully. See  
 1537           descriptions of these in the Base Definitions volume of IEEE Std. 1003.1-200x, Section 6.2, Character  
 1538           Encoding. If such character encodings are supported, it is expected that this will impact the Base  
 1539           Definitions volume of IEEE Std. 1003.1-200x, Section 6.2, Character Encoding, the Base Definitions  
 1540           volume of IEEE Std. 1003.1-200x, Chapter 7, Locale, the Base Definitions volume of  
 1541           IEEE Std. 1003.1-200x, Chapter 9, Regular Expressions, and the comm, cut, diff, grep, head, join,  
 1542           paste, and tail utilities.

1543           The character set description file provides:

- 1544           • The capability to describe character set attributes (such as collation order or character  
 1545           classes) independent of character set encoding, and using only the characters in the portable  
 1546           character set. This makes it possible to create generic *localedef* source files for all codesets that  
 1547           share the portable character set (such as the ISO 8859 family or IBM Extended ASCII).
- 1548           • Standardized symbolic names for all characters in the portable character set, making it  
 1549           possible to refer to any such character regardless of encoding.

1550           Implementations are free to choose their own symbolic names, as long as the names identified  
 1551           by this volume of IEEE Std. 1003.1-200x are also defined; this provides support for already  
 1552           existing “character names”.

1553           The names selected for the members of the portable character set follow the  
 1554           ISO/IEC 8859-1:1998 standard and the ISO/IEC 10646-1:1993 standard. However, several  
 1555           commonly used UNIX system names occur as synonyms in the list:

- 1556           • The historical UNIX system names are used for control characters.
- 1557           • The word “slash” is given in addition to “solidus”.
- 1558           • The word “backslash” is given in addition to “reverse-solidus”.
- 1559           • The word “hyphen” is given in addition to “hyphen-minus”.
- 1560           • The word “period” is given in addition to “full-stop”.
- 1561           • For digits, the word “digit” is eliminated.
- 1562           • For letters, the words “Latin Capital Letter” and “Latin Small Letter” are eliminated.
- 1563           • The words “left brace” and “right brace” are given in addition to “left-curly-bracket” and  
 1564           “right-curly-bracket”.
- 1565           • The names of the digits are preferred over the numbers to avoid possible confusion between  
 1566           ‘0’ and ‘o’, and between ‘1’ and ‘l’ (one and the letter ell).

1567           The names for the control characters in the Base Definitions volume of IEEE Std. 1003.1-200x,  
 1568           Chapter 6, Character Set were taken from the ISO/IEC 4873:1991 standard.

1569           The charmap file was introduced to resolve problems with the portability of, especially, *localedef*  
 1570           sources. IEEE Std. 1003.1-200x assumes that the portable character set is constant across all  
 1571           locales, but does not prohibit implementations from supporting two incompatible codings, such  
 1572           as both ASCII and EBCDIC. Such dual-support implementations should have all charmaps and  
 1573           *localedef* sources encoded using one portable character set, in effect cross-compiling for the other  
 1574           environment. Naturally, charmaps (and *localedef* sources) are only portable without  
 1575           transformation between systems using the same encodings for the portable character set. They  
 1576           can, however, be transformed between two sets using only a subset of the actual characters (the  
 1577           portable character set). However, the particular coded character set used for an application or an  
 1578           implementation does not necessarily imply different characteristics or collation; on the contrary,  
 1579           these attributes should in many cases be identical, regardless of codeset. The charmap provides

1580 the capability to define a common locale definition for multiple codesets (the same *localedef*  
 1581 source can be used for codesets with different extended characters; the ability in the charmap to  
 1582 define empty names allows for characters missing in certain codesets).

1583 The **<escape\_char>** declaration was added at the request of the international community to ease  
 1584 the creation of portable charmap files on terminals not implementing the default backslash  
 1585 escape. The **<comment\_char>** declaration was added at the request of the international  
 1586 community to eliminate the potential confusion between the number sign and the pound sign.

1587 The octal number notation with no leading zero required was selected to match those of *awk* and  
 1588 *tr* and is consistent with that used by *localedef*. To avoid confusion between an octal constant  
 1589 and the back-references used in *localedef* source, the octal, hexadecimal, and decimal constants  
 1590 shall contain at least two digits. As single-digit constants are relatively rare, this should not  
 1591 impose any significant hardship. Provision is made for more digits to account for systems in  
 1592 which the byte size is larger than 8 bits. For example, a Unicode (ISO/IEC 10646-1:1993  
 1593 standard) system that has defined 16-bit bytes may require six octal, four hexadecimal, and five  
 1594 decimal digits.

1595 The decimal notation is supported because some newer international standards define character  
 1596 values in decimal, rather than in the old column/row notation.

1597 The charmap identifies the coded character sets supported by an implementation. At least one  
 1598 charmap shall be provided, but no implementation is required to provide more than one.  
 1599 Likewise, implementations can allow users to generate new charmaps (for instance, for a new  
 1600 version of the ISO 8859 family of coded character sets), but does not have to do so. If users are  
 1601 allowed to create new charmaps, the system documentation describes the rules that apply (for  
 1602 instance, “only coded character sets that are supersets of the ISO/IEC 646:1991 standard IRV, no  
 1603 multi-byte characters”).

1604 This addition of the **WIDTH** specification satisfies the following requirement from the  
 1605 ISO POSIX-2:1993 standard, Annex H.1:

1606 (9) *The definition of column position relies on the implementation’s knowledge of the integral width*  
 1607 *of the characters. The charmap or LC\_CTYPE locale definitions should be enhanced to allow*  
 1608 *application specification of these widths.*

1609 The character “width” information was first considered for inclusion under *LC\_CTYPE* but was  
 1610 moved because it is more closely associated with the information in the *charmap* than  
 1611 information in the locale source (cultural conventions information). Concerns were raised that  
 1612 formalizing this type of information is moving the locale source definition from the codeset-  
 1613 independent entity that it was designed to be to a repository of codeset-specific information. A  
 1614 similar issue occurred with the **<code\_set\_name>**, **<mb\_cur\_max>**, and **<mb\_cur\_min>**  
 1615 information, which was resolved to reside in the *charmap* definition.

1616 The width definition was added to the IEEE P1003.2b draft standard with the intent that the  
 1617 *wcswidth()* and/or *wcwidth()* functions (currently specified in the System Interfaces volume of  
 1618 IEEE Std. 1003.1-200x) be the mechanism to retrieve the character width information.

## 1619 A.7 Locale

### 1620 A.7.1 General

1621 The description of locales is based on work performed in the UniForum Technical Committee  
1622 Subcommittee on Internationalization. Wherever appropriate, keywords are taken from the  
1623 ISO C standard or the X/Open Portability Guide.

1624 The value used to specify a locale with environment variables is the name specified as the *name*  
1625 operand to the *localedef* utility when the locale was created. This provides a verifiable method to  
1626 create and invoke a locale.

1627 The “object” definitions need not be portable, as long as “source” definitions are. Strictly  
1628 speaking, source definitions are portable only between implementations using the same  
1629 character set(s). Such source definitions, if they use symbolic names only, easily can be ported  
1630 between systems using different codesets, as long as the characters in the portable character set  
1631 (see the Base Definitions volume of IEEE Std. 1003.1-200x, Section 6.1, Portable Character Set )  
1632 have common values between the codesets; this is frequently the case in historical  
1633 implementations. Of source, this requires that the symbolic names used for characters outside  
1634 the portable character set be identical between character sets. The definition of symbolic names  
1635 for characters is outside the scope of IEEE Std. 1003.1-200x, but is certainly within the scope of  
1636 other standards organizations.

1637 Applications can select the desired locale by invoking the *setlocale()* function (or equivalent)  
1638 with the appropriate value. If the function is invoked with an empty string, the value of the  
1639 corresponding environment variable is used. If the environment variable is not set or is set to the  
1640 empty string, the implementation sets the appropriate environment as defined in the Base  
1641 Definitions volume of IEEE Std. 1003.1-200x, Chapter 8, Environment Variables.

### 1642 A.7.2 POSIX Locale

1643 The POSIX locale is equal to the C locale. To avoid being classified as a C-language function, the  
1644 name has been changed to the POSIX locale; the environment variable value can be either  
1645 "POSIX" or, for historical reasons, "C".

1646 The POSIX definitions mirror the historical UNIX system behavior.

1647 The use of symbolic names for characters in the tables does not imply that the POSIX locale must  
1648 be described using symbolic character names, but merely that it may be advantageous to do so.

### 1649 A.7.3 Locale Definition

1650 The decision to separate the file format from the *localedef* utility description was only partially  
1651 editorial. Implementations may provide other interfaces than *localedef*. Requirements on “the  
1652 utility”, mostly concerning error messages, are described in this way because they are meant to  
1653 affect the other interfaces implementations may provide as well as *localedef*.

1654 The text about POSIX2\_LOCALEDEF does not mean that internationalization is optional; only  
1655 that the functionality of the *localedef* utility is. REs, for instance, must still be able to recognize,  
1656 for example, character class expressions such as "[[:alpha:]]". A possible analogy is with  
1657 an applications development environment; while all conforming implementations must be  
1658 capable of executing applications, not all need to have the development environment installed.  
1659 The assumption is that the capability to modify the behavior of utilities (and applications) via  
1660 locale settings must be supported. If the *localedef* utility is not present, then the only choice is to  
1661 select an existing (presumably implementation-documented) locale. An implementation could,  
1662 for example, choose to support only the POSIX locale, which would in effect limit the amount of

1663 changes from historical implementations quite drastically. The *localedef* utility is still required,  
1664 but would always terminate with an exit code indicating that no locale could be created.  
1665 Supported locales must be documented using the syntax defined in this chapter. (This ensures  
1666 that users can accurately determine what capabilities are provided. If the implementation  
1667 decides to provide additional capabilities to the ones in this chapter, that is already provided  
1668 for.)

1669 If the option is present (that is, locales can be created), then the *localedef* utility must be capable  
1670 of creating locales based on the syntax and rules defined in this chapter. This does not mean that  
1671 the implementation cannot also provide alternate means for creating locales.

1672 The octal, decimal, and hexadecimal notations are the same employed by the charmap facility  
1673 (see the Base Definitions volume of IEEE Std. 1003.1-200x, Section 6.4, Character Set Description  
1674 File). To avoid confusion between an octal constant and a back-reference, the octal, hexadecimal,  
1675 and decimal constants must contain at least two digits. As single-digit constants are relatively  
1676 rare, this should not impose any significant hardship. Provision is made for more digits to  
1677 account for systems in which the byte size is larger than 8 bits. For example, a Unicode (see the  
1678 ISO/IEC 10646-1:1993 standard) system that has defined 16-bit bytes may require six octal, four  
1679 hexadecimal, and five decimal digits. As with the charmap file, multi-byte characters are  
1680 described in the locale definition file using “big-endian” notation for reasons of portability.  
1681 There is no requirement that the internal representation in the computer memory be in this same  
1682 order.

1683 One of the guidelines used for the development of this volume of IEEE Std. 1003.1-200x is that  
1684 characters outside the invariant part of the ISO/IEC 646:1991 standard should not be used in  
1685 portable specifications. The backslash character is not in the invariant part; the number sign is,  
1686 but with multiple representations: as a number sign, and as a pound sign. As far as general  
1687 usage of these symbols, they are covered by the “grandfather clause”, but for newly defined  
1688 interfaces, the WG15 POSIX working group has requested that POSIX provide alternate  
1689 representations. Consequently, while the default escape character remains the backslash and the  
1690 default comment character is the number sign, implementations are required to recognize  
1691 alternative representations, identified in the applicable source file via the `<escape_char>` and  
1692 `<comment_char>` keywords.

### 1693 A.7.3.1 *LC\_CTYPE*

1694 The *LC\_CTYPE* category is primarily used to define the encoding-independent aspects of a  
1695 character set, such as character classification. In addition, certain encoding-dependent  
1696 characteristics are also defined for an application via the *LC\_CTYPE* category.  
1697 IEEE Std. 1003.1-200x does not mandate that the encoding used in the locale is the same as the  
1698 one used by the application because an implementation may decide that it is advantageous to  
1699 define locales in a system-wide encoding rather than having multiple, logically identical locales  
1700 in different encodings, and to convert from the application encoding to the system-wide  
1701 encoding on usage. Other implementations could require encoding-dependent locales.

1702 In either case, the *LC\_CTYPE* attributes that are directly dependent on the encoding, such as  
1703 `<mb_cur_max>` and the display width of characters, are not user-specifiable in a locale source  
1704 and are consequently not defined as keywords.

1705 Implementations may define additional keywords or extend the *LC\_CTYPE* mechanism to allow  
1706 application-defined keywords.

1707 The text “The ellipsis specification shall only be valid within a single encoded character set” is  
1708 present because it is possible to have a locale supported by multiple character encodings, as  
1709 explained in the rationale for the Base Definitions volume of IEEE Std. 1003.1-200x, Section 6.1,  
1710 Portable Character Set. An example given there is of a possible Japanese-based locale supported

1711 by a mixture of the character sets JIS X 0201 Roman, JIS X 0208, and JIS X 0201 Katakana.  
 1712 Attempting to express a range of characters across these sets is not logical and the  
 1713 implementation is free to reject such attempts.

1714 As the `LC_CTYPE` character classes are based on the ISO C standard character class definition,  
 1715 the category does not support multi-character elements. For instance, the German character  
 1716 <sharp-s> is traditionally classified as a lowercase letter. There is no corresponding uppercase  
 1717 letter; in proper capitalization of German text, the <sharp-s> will be replaced by "SS"; that is, by  
 1718 two characters. This kind of conversion is outside the scope of the **toupper** and **tolower**  
 1719 keywords.

1720 Where IEEE Std. 1003.1-200x specifies that only certain characters can be specified, as for the  
 1721 keywords **digit** and **xdigit**, the specified characters shall be from the portable character set, as  
 1722 shown. As an example, only the Arabic digits 0 through 9 are acceptable as digits.

1723 The character classes **digit**, **xdigit**, **lower**, **upper**, and **space** have a set of automatically included  
 1724 characters. These only need to be specified if the character values (that is, encoding) differs from  
 1725 the implementation default values. It is not possible to define a locale without these  
 1726 automatically included characters unless some implementation extension is used to prevent  
 1727 their inclusion. Such a definition would not be a proper superset of the C locale, and thus, it  
 1728 might not be possible for the standard utilities to be implemented as programs conforming to  
 1729 the ISO C standard.

1730 The definition of character class **digit** requires that only ten characters—the ones defining  
 1731 digits—can be specified; alternate digits (for example, Hindi or Kanji) cannot be specified here.  
 1732 However, the encoding may vary if an implementation supports more than one encoding.

1733 The definition of character class **xdigit** requires that the characters included in character class  
 1734 **digit** are included here also and allows for different symbols for the hexadecimal digits 10  
 1735 through 15.

1736 The inclusion of the **charclass** keyword satisfies the following requirement from the  
 1737 ISO POSIX-2: 1993 standard, Annex H.1:

1738 (3) *The `LC_CTYPE` (2.5.2.1) locale definition should be enhanced to allow user-specified additional*  
 1739 *character classes, similar in concept to the ISO C standard Multibyte Support Extension (MSE)*  
 1740 *is\_wctype() function.*

1741 This keyword was previously included in The Open Group specifications and is now mandated  
 1742 in the Shell and Utilities volume of IEEE Std. 1003.1-200x.

1743 The symbolic constant `{CHARCLASS_NAME_MAX}` was also adopted from The Open Group  
 1744 specifications. Application portability is enhanced by the use of symbolic constants.

#### 1745 A.7.3.2 `LC_COLLATE`

1746 The rules governing collation depend to some extent on the use. At least five different levels of  
 1747 increasingly complex collation rules can be distinguished:

- 1748 1. *Byte/machine code order*: This is the historical collation order in the UNIX system and many  
 1749 proprietary operating systems. Collation is here performed character by character, without  
 1750 any regard to context. The primary virtue is that it usually is quite fast and also  
 1751 completely deterministic; it works well when the native machine collation sequence  
 1752 matches the user expectations.
- 1753 2. *Character order*: On this level, collation is also performed character by character, without  
 1754 regard to context. The order between characters is, however, not determined by the code  
 1755 values, but on the expectations by the user of the “correct” order between characters. In

1756 addition, such a (simple) collation order can specify that certain characters collate equally  
1757 (for example, uppercase and lowercase letters).

1758 3. *String ordering*: On this level, entire strings are compared based on relatively  
1759 straightforward rules. Several “passes” may be required to determine the order between  
1760 two strings. Characters may be ignored in some passes, but not in others; the strings may  
1761 be compared in different directions; and simple string substitutions may be performed  
1762 before strings are compared. This level is best described as “dictionary” ordering; it is  
1763 based on the spelling, not the pronunciation, or meaning, of the words.

1764 4. *Text search ordering*: This is a further refinement of the previous level, best described as  
1765 “telephone book ordering”; some common homonyms (words spelled differently but with  
1766 the same pronunciation) are collated together; numbers are collated as if they were spelled  
1767 out, and so on.

1768 5. *Semantic-level ordering*: Words and strings are collated based on their meaning; entire words  
1769 (such as “the”) are eliminated; the ordering is not deterministic. This usually requires  
1770 special software and is highly dependent on the intended use.

1771 While the historical collation order formally is at level 1, for the English language it corresponds  
1772 roughly to elements at level 2. The user expects to see the output from the *ls* utility sorted very  
1773 much as it would be in a dictionary. While telephone book ordering would be an optimal goal  
1774 for standard collation, this was ruled out as the order would be language-dependent.  
1775 Furthermore, a requirement was that the order must be determined solely from the text string  
1776 and the collation rules; no external information (for example, “pronunciation dictionaries”)   
1777 could be required.

1778 As a result, the goal for the collation support is at level 3. This also matches the requirements for  
1779 the Canadian collation order, as well as other, known collation requirements for alphabetic  
1780 scripts. It specifically rules out collation based on pronunciation rules or based on semantic  
1781 analysis of the text.

1782 The syntax for the *LC\_COLLATE* category source meets the requirements for level 3 and has  
1783 been verified to produce the correct result with examples based on French, Canadian, and  
1784 Danish collation order. Because it supports multi-character collating elements, it is also capable  
1785 of supporting collation in codesets where a character is expressed using non-spacing characters  
1786 followed by the base character (such as the ISO/IEC 6937: 1994 standard).

1787 The directives that can be specified in an operand to the **order\_start** keyword are based on the  
1788 requirements specified in several proposed standards and in customary use. The following is a  
1789 rephrasing of rules defined for “lexical ordering in English and French” by the Canadian  
1790 Standards Association (the text in square brackets is rephrased):

1791 • Once special characters [punctuation] have been removed from original strings, the ordering  
1792 is determined by scanning forwards (left to right) [disregarding case and diacriticals].

1793 • In case of equivalence, special characters are once again removed from original strings and  
1794 the ordering is determined by scanning backwards (starting from the rightmost character of  
1795 the string and back), character by character [disregarding case but considering diacriticals].

1796 • In case of repeated equivalence, special characters are removed again from original strings  
1797 and the ordering is determined by scanning forwards, character by character [considering  
1798 both case and diacriticals].

1799 • If there is still an ordering equivalence after the first three rules have been applied, then only  
1800 special characters and the position they occupy in the string are considered to determine  
1801 ordering. The string that has a special character in the lowest position comes first. If two  
1802 strings have a special character in the same position, the character [with the lowest collation



1803 value] comes first. In case of equality, the other special characters are considered until there  
1804 is a difference or until all special characters have been exhausted.

1805 It is estimated that this part of IEEE Std. 1003.1-200x covers the requirements for all European  
1806 languages, and no particular problems are anticipated with Slavic or Middle East character sets.

1807 The Far East (particularly Japanese/Chinese) collations are often based on contextual  
1808 information and pronunciation rules (the same ideogram can have different meanings and  
1809 different pronunciations). Such collation, in general, falls outside the desired goal of  
1810 IEEE Std. 1003.1-200x. There are, however, several other collation rules (stroke/radical or “most  
1811 common pronunciation”) that can be supported with the mechanism described here.

1812 The character (and collating element) order is defined by the order in which characters and  
1813 elements are specified between the **order\_start** and **order\_end** keywords. This character order is  
1814 used in range expressions in REs (see the Base Definitions volume of IEEE Std. 1003.1-200x,  
1815 Chapter 9, Regular Expressions). Weights assigned to the characters and elements define the  
1816 collation sequence; in the absence of weights, the character order is also the collation sequence.

#### 1817 A.7.3.3 *LC\_MONETARY*

1818 The currency symbol does not appear in *LC\_MONETARY* because it is not defined in the C locale  
1819 of the ISO C standard.

1820 The ISO C standard limits the size of decimal points and thousands delimiters to single-byte  
1821 values. In locales based on multi-byte coded character sets, this cannot be enforced;  
1822 IEEE Std. 1003.1-200x does not prohibit such characters, but makes the behavior unspecified (in  
1823 the text “In contexts where other standards ...”).

1824 The grouping specification is based on, but not identical to, the ISO C standard. The -1 signals  
1825 that no further grouping shall be performed; the equivalent of {CHAR\_MAX} in the ISO C  
1826 standard.

1827 The text “the value is not available in the locale” is taken from the ISO C standard and is used  
1828 instead of the “unspecified” text in early proposals. There is no implication that omitting these  
1829 keywords or assigning them values of " " or -1 produces unspecified results; such omissions or  
1830 assignments eliminate the effects described for the keyword or produce zero-length strings, as  
1831 appropriate.

1832 The locale definition is an extension of the ISO C standard *localeconv()* specification. In  
1833 particular, rules on how **currency\_symbol** is treated are extended to also cover **int\_curr\_symbol**,  
1834 and **p\_set\_by\_space** and **n\_sep\_by\_space** have been augmented with the value 2, which places  
1835 a <space> between the sign and the symbol (if they are adjacent; otherwise, it should be treated  
1836 as a 0).

#### 1837 A.7.3.4 *LC\_NUMERIC*

1838 See the rationale for *LC\_MONETARY* for a description of the behavior of grouping.

#### 1839 A.7.3.5 *LC\_TIME*

1840 Although certain of the field descriptors in the POSIX locale (such as the name of the month) are  
1841 shown with initial capital letters, this need not be the case in other locales. Programs using these  
1842 fields may need to adjust the capitalization if the output is going to be used at the beginning of a  
1843 sentence.

1844 The *LC\_TIME* descriptions of **abday**, **day**, **mon**, and **abmon** imply a Gregorian style calendar (7-  
1845 day weeks, 12-month years, leap years, and so on). Formatting time strings for other types of  
1846 calendars is outside the scope of IEEE Std. 1003.1-200x.

1847 While the ISO 8601:1988 standard numbers the weekdays starting with Monday, historical  
 1848 practice is to use the Sunday as the first day. Rather than change the order and introduce  
 1849 potential confusion, the days must be specified beginning with Sunday; previous references to  
 1850 “first day” have been removed. Note also that the Shell and Utilities volume of  
 1851 IEEE Std. 1003.1-200x *date* utility supports numbering compliant with the ISO 8601:1988  
 1852 standard.

1853 As specified under *date* in the Shell and Utilities volume of IEEE Std. 1003.1-200x and *strftime()*  
 1854 in the System Interfaces volume of IEEE Std. 1003.1-200x, the field descriptors corresponding to  
 1855 the optional keywords consist of a modifier followed by a traditional field descriptor (for  
 1856 instance %Ex). If the optional keywords are not supported by the implementation or are  
 1857 unspecified for the current locale, these field descriptors are treated as the traditional field  
 1858 descriptor. For example, assume the following keywords:

```
1859 alt_digits "0th";"1st";"2nd";"3rd";"4th";"5th";\  
1860           "6th";"7th";"8th";"9th";"10th"
```

```
1861 d_fmt      "The %Od day of %B in %Y"
```

1862 On 7/4/1776, the %x field descriptor would result in "The 4th day of July in 1776",  
 1863 while on 7/14/1789 would result in "The 14 day of July in 1789". It can be noted that  
 1864 the above example is for illustrative purposes only; the %O modifier is primarily intended to  
 1865 provide for Kanji or Hindi digits in *date* formats.

1866 A.7.3.6 *LC\_MESSAGES*

#### 1867 **A.7.4 Locale Definition Grammar**

1868 There is no additional rationale for this section.

1869 A.7.4.1 *Locale Lexical Conventions*

1870 There is no additional rationale for this section.

1871 A.7.4.2 *Locale Grammar*

1872 There is no additional rationale for this section.

#### 1873 **A.7.5 Locale Definition Example**

1874 There is no additional rationale for this section.

1875 **A.8 Environment Variables**1876 **A.8.1 Environment Variable Definition**

1877 The variable *environ* is not intended to be declared in any header, but rather to be declared by the  
 1878 user for accessing the array of strings that is the environment. This is the traditional usage of the  
 1879 symbol. Putting it into a header could break some programs that use the symbol for their own  
 1880 purposes.

1881 The decision to restrict conforming systems to the use of digits, uppercase letters, and  
 1882 underscores for environment variable names allows applications to use lowercase letters in their  
 1883 environment variable names without conflicting with any conforming system.

1884 **A.8.2 Internationalization Variables**

1885 The text about locale implies that any utilities written in standard C and conforming to  
 1886 IEEE Std. 1003.1-200x must issue the following call:

```
1887     setlocale(LC_ALL, "")
```

1888 If this were omitted, the ISO C standard specifies that the C locale would be used.

1889 If any of the environment variables are invalid, it makes sense to default to an implementation-  
 1890 defined, consistent locale environment. It is more confusing for a user to have partial settings  
 1891 occur in case of a mistake. All utilities would then behave in one language/cultural  
 1892 environment. Furthermore, it provides a way of forcing the whole environment to be the  
 1893 implementation-defined default. Disastrous results could occur if a pipeline of utilities partially  
 1894 uses the environment variables in different ways. In this case, it would be appropriate for  
 1895 utilities that use *LANG* and related variables to exit with an error if any of the variables are  
 1896 invalid. For example, users typing individual commands at a terminal might want *date* to work if  
 1897 *LC\_MONETARY* is invalid as long as *LC\_TIME* is valid. Since these are conflicting reasonable  
 1898 alternatives, IEEE Std. 1003.1-200x leaves the results unspecified if the locale environment  
 1899 variables would not produce a complete locale matching the specification of the user.

1900 The locale settings of individual categories cannot be truly independent and still guarantee  
 1901 correct results. For example, when collating two strings, characters must first be extracted from  
 1902 each string (governed by *LC\_CTYPE*) before being mapped to collating elements (governed by  
 1903 *LC\_COLLATE*) for comparison. That is, if *LC\_CTYPE* is causing parsing according to the rules of  
 1904 a large, multi-byte code set (potentially returning 20 000 or more distinct character codeset  
 1905 values), but *LC\_COLLATE* is set to handle only an 8-bit codeset with 256 distinct characters,  
 1906 meaningful results are obviously impossible.

1907 The *LC\_MESSAGES* variable affects the language of messages generated by the standard  
 1908 utilities.

1909 The description of the environment variable names starting with the characters “LC\_”  
 1910 acknowledges the fact that the interfaces presented may be extended as new international  
 1911 functionality is required. In the ISO C standard, names preceded by “LC\_” are reserved in the  
 1912 name space for future categories.

1913 To avoid name clashes, new categories and environment variables are divided into two  
 1914 classifications: *implementation-independent* and *implementation-defined*.

1915 Implementation-independent names will have the following format:

```
1916     LC_NAME
```

1917 where *NAME* is the name of the new category and environment variable. Capital letters must be  
 1918 used for implementation-independent names.

1919 Implementation-defined names must be in lowercase letters, as below:

1920 `LC_name`

### 1921 **A.8.3 Other Environment Variables**

1922 The quoted form of the timezone variable allows timezone names of the form UTC+1 (or any  
 1923 name that contains the character plus ('+'), the character minus ('-'), or digits), which may be  
 1924 appropriate for countries that do not have an official timezone name. It would be coded as  
 1925 <UTC+1>+1<UTC+2>, which would cause *std* to have a value of UTC+1 and *dst* a value of  
 1926 UTC+2, each with a length of 6 characters. This does not appear to conflict with any existing  
 1927 usage. The characters '<' and '>' were chosen for quoting because they are easier to parse  
 1928 visually than a quoting character that does not provide some sense of bracketing (and in a string  
 1929 like this, such bracketing is helpful). They were also chosen because they do not need special  
 1930 treatment when assigning to the *TZ* variable. Users are often confused by embedding quotes in a  
 1931 string. Because '<' and '>' are meaningful to the shell, the whole string would have to be  
 1932 quoted, but that is easily explained. (Parentheses would have presented the same problems.)  
 1933 Although the '>' symbol could have been permitted in the string by either escaping it or  
 1934 doubling it, it seemed of little value to require that. This could be provided as an extension if  
 1935 there was a need. Timezone names of this new form lead to a requirement that the value of  
 1936 `{_POSIX_TZNAME_MAX}` change from 3 to 6.

### 1937 **COLUMNS, LINES**

1938 The default value for the number of column positions, *COLUMNS*, and screen height, *LINES*, are  
 1939 unspecified because historical implementations use different methods to determine values  
 1940 corresponding to the size of the screen in which the utility is run. This size is typically known to  
 1941 the implementation through the value of *TERM*, or by more elaborate methods such as  
 1942 extensions to the *stty* utility or knowledge of how the user is dynamically resizing windows on a  
 1943 bit-mapped display terminal. Users should not need to set these variables in the environment  
 1944 unless there is a specific reason to override the default behavior of the implementation, such as  
 1945 to display data in an area arbitrarily smaller than the terminal or window. Values for these  
 1946 variables that are not decimal integers greater than zero are implicitly undefined values; it is  
 1947 unnecessary to enumerate all of the possible values outside of the acceptable set.

### 1948 **PATH**

1949 Many historical implementations of the Bourne shell do not interpret a trailing colon to represent  
 1950 the current working directory and are thus non-conforming. The C Shell and the KornShell  
 1951 conform to IEEE Std. 1003.1-200x on this point. The usual name of dot may also be used to refer  
 1952 to the current working directory.

1953 Many implementations historically have used a default value of `/bin` and `/usr/bin` for the *PATH*  
 1954 variable. IEEE Std. 1003.1-200x does not mandate this default path be identical to that retrieved  
 1955 from `getconf _CS_PATH` because it is likely that the standardized utilities may be provided in  
 1956 another directory separate from the directories used by some historical applications.

1957 **LOGNAME**

1958 In most implementations, the value of such a variable is easily forged, so security-critical  
1959 applications should rely on other means of determining user identity. *LOGNAME* is required to  
1960 be constructed from the portable file name character set for reasons of interchange. No  
1961 diagnostic condition is specified for violating this rule, and no requirement for enforcement  
1962 exists. The intent of the requirement is that if extended characters are used, the “guarantee” of  
1963 portability implied by a standard is void.

1964 **SHELL**

1965 The *SHELL* variable names the preferred shell of the user; it is a guide to applications. There is  
1966 no direct requirement that that shell conform to IEEE Std. 1003.1-200x; that decision should rest  
1967 with the user. It is the intention of the standard developers that alternative shells be permitted, if  
1968 the user chooses to develop or acquire one. An operating system that builds its shell into the  
1969 “kernel” in such a manner that alternative shells would be impossible does not conform to the  
1970 spirit of IEEE Std. 1003.1-200x.

1971 **CHANGE HISTORY**1972 **Issue 6**

1973 Changed format of *TZ* field to allow for the quoted form as defined in previous  
1974 versions of the ISO POSIX-1 standard.

1975 **A.9 Regular Expressions**

1976 Rather than repeating the description of REs for each utility supporting REs, the standard  
 1977 developers preferred a common, comprehensive description of regular expressions in one place.  
 1978 The most common behavior is described here, and exceptions or extensions to this are  
 1979 documented for the respective utilities, as appropriate.

1980 The BRE corresponds to the *ed* or historical *grep* type, and the ERE corresponds to the historical  
 1981 *egrep* type (now *grep -E*).

1982 The text is based on the *ed* description and substantially modified, primarily to aid developers  
 1983 and others in the understanding of the capabilities and limitations of REs. Much of this was  
 1984 influenced by internationalization requirements.

1985 It should be noted that the definitions in this section do not cover the *tr* utility; the *tr* syntax does  
 1986 not employ REs.

1987 The specification of REs is particularly important to internationalization because pattern  
 1988 matching operations are very basic operations in business and other operations. The syntax and  
 1989 rules of REs are intended to be as intuitive as possible to make them easy to understand and use.  
 1990 The historical rules and behavior do not provide that capability to non-English language users,  
 1991 and do not provide the necessary support for commonly used characters and language  
 1992 constructs. It was necessary to provide extensions to the historical RE syntax and rules to  
 1993 accommodate other languages.

1994 As they are limited to bracket expressions, the rationale for these modifications is in the Base  
 1995 Definitions volume of IEEE Std. 1003.1-200x, Section 9.3.5, RE Bracket Expression.

1996 **A.9.1 Regular Expression Definitions**

1997 It is possible to determine what strings correspond to subexpressions by recursively applying  
 1998 the leftmost longest rule to each subexpression, but only with the proviso that the overall match  
 1999 is leftmost longest. For example, matching "`\(ac*\\)c*d[ac]*\1`" against *acdacaaa* matches  
 2000 *acdacaaa* (with `\1=a`); simply matching the longest match for "`\(ac*\\)`" would yield `\1=ac`, but  
 2001 the overall match would be smaller (*acdac*). Conceptually, the implementation must examine  
 2002 every possible match and among those that yield the leftmost longest total matches, pick the one  
 2003 that does the longest match for the leftmost subexpression, and so on. Note that this means that  
 2004 matching by subexpressions is context-dependent: a subexpression within a larger RE may  
 2005 match a different string from the one it would match as an independent RE, and two instances of  
 2006 the same subexpression within the same larger RE may match different lengths even in similar  
 2007 sequences of characters. For example, in the ERE "`(a.*b)(a.*b)`", the two identical  
 2008 subexpressions would match four and six characters, respectively, of *accbacccb*.

2009 The definition of *single character* has been expanded to include also collating elements consisting  
 2010 of two or more characters; this expansion is applicable only when a bracket expression is  
 2011 included in the BRE or ERE. An example of such a collating element may be the Dutch *ij*, which  
 2012 collates as a 'y'. In some encodings, a ligature "i with j" exists as a character and would  
 2013 represent a single-character collating element. In another encoding, no such ligature exists, and  
 2014 the two-character sequence *ij* is defined as a multi-character collating element. Outside brackets,  
 2015 the *ij* is treated as a two-character RE and matches the same characters in a string. Historically, a  
 2016 bracket expression only matched a single character. If, however, the bracket expression defines,  
 2017 for example, a range that includes *ij*, then this particular bracket expression also matches a  
 2018 sequence of the two characters 'i' and 'j' in the string.

2019 **A.9.2 Regular Expression General Requirements**

2020 The definition of which sequence is matched when several are possible is based on the leftmost-  
 2021 longest rule historically used by deterministic recognizers. This rule is easier to define and  
 2022 describe, and arguably more useful, than the first-match rule historically used by non-  
 2023 deterministic recognizers. It is thought that dependencies on the choice of rule are rare; carefully  
 2024 contrived examples are needed to demonstrate the difference.

2025 A formal expression of the leftmost-longest rule is:

2026       The search is performed as if all possible suffixes of the string were tested for a prefix  
 2027       matching the pattern; the longest suffix containing a matching prefix is chosen, and the  
 2028       longest possible matching prefix of the chosen suffix is identified as the matching sequence.

2029 Historically, most RE implementations only match lines, not strings. However, that is more an  
 2030 effect of the usage than of an inherent feature of REs themselves. Consequently, IEEE Std. 1003.1-200x  
 2031 does not regard <newline>s as special; they are ordinary characters, and both a period and a non-matching  
 2032 list can match them. Those utilities (like *grep*) that do not allow <newline>s to match are responsible  
 2033 for eliminating any <newline> from strings before matching against the RE. The *regcomp()* function,  
 2034 however, can provide support for such processing without violating the rules of this section.  
 2035

2036 The definition of case-insensitive processing is intended to allow matching of multi-character  
 2037 collating elements as well as characters. For instance, as each character in the string is matched  
 2038 using both its cases, the RE "[ [ . Ch . ] ]", when matched against "char", is in reality matched  
 2039 against "ch", "Ch", "cH", and "CH".

2040 Some implementations of *egrep* have had very limited flexibility in handling complex EREs. IEEE Std.  
 2041 1003.1-200x does not attempt to define the complexity of a BRE or ERE, but does place a lower limit  
 2042 on it—any RE must be handled, as long as it can be expressed in 256 bytes or less. (Of course, this  
 2043 does not place an upper limit on the implementation.) There are historical programs using a non-  
 2044 deterministic-recognizer implementation that should have no difficulty with this limit. It is possible  
 2045 that a good approach would be to attempt to use the faster, but more limited, deterministic  
 2046 recognizer for simple expressions and to fall back on the non-deterministic recognizer for those  
 2047 expressions requiring it. Non-deterministic implementations must be careful to observe the rules on  
 2048 which match is chosen; the longest match, not the first match, starting at a given character is used.  
 2049

2050 The term *invalid* highlights a difference between this section and some others: IEEE Std. 1003.1-200x  
 2051 frequently avoids mandating of errors for syntax violations because they can be used by implementors  
 2052 to trigger extensions. However, the authors of the internationalization features of REs wanted to  
 2053 mandate errors for certain conditions to identify usage problems or non-portable constructs. These  
 2054 are identified within this rationale as appropriate. The remaining syntax violations have been left  
 2055 implicitly or explicitly undefined. For example, the BRE construct "\{1,2,3\}" does not comply with  
 2056 the grammar. A conforming application cannot rely on it producing an error nor matching the literal  
 2057 characters "\{1,2,3\}". The term “undefined” was used in favor of “unspecified” because many of  
 2058 the situations are considered errors on some implementations, and the standard developers  
 2059 considered that consistency throughout the section was preferable to mixing undefined and  
 2060 unspecified.  
 2061

2062 **A.9.3 Basic Regular Expressions**

2063 There is no additional rationale for this section.

2064 **A.9.3.1 BREs Matching a Single Character or Collating Element**

2065 There is no additional rationale for this section.

2066 **A.9.3.2 BRE Ordinary Characters**

2067 There is no additional rationale for this section.

2068 **A.9.3.3 BRE Special Characters**

2069 There is no additional rationale for this section.

2070 **A.9.3.4 Periods in BREs**

2071 There is no additional rationale for this section.

2072 **A.9.3.5 RE Bracket Expression**

2073 Range expressions are, historically, an integral part of REs. However, the requirements of  
 2074 “natural language behavior” and portability do conflict: ranges must be treated according to the  
 2075 current collating sequence and include such characters that fall within the range based on that  
 2076 collating sequence, regardless of character values. This means, however, that the interpretation  
 2077 will differ depending on collating sequence. If, for instance, one collating sequence defines ‘a’ as  
 2078 a variant of ‘a’, while another defines it as a letter following ‘z’, then the expression “[a-z]”  
 2079 is valid in the first language and invalid in the second. This kind of ambiguity should be avoided  
 2080 in portable applications, and therefore the standard developers elected to state that ranges must  
 2081 not be used in strictly conforming applications; however, implementations must support them.

2082 Some historical implementations allow range expressions where the ending range point of one  
 2083 range is also the starting point of the next (for instance, “[a-m-o]”). This behavior should not  
 2084 be permitted, but to avoid breaking historical implementations, it is now *undefined* whether it is a  
 2085 valid expression and how it should be interpreted.

2086 Current practice in *awk* and *lex* is to accept escape sequences in bracket expressions as per the  
 2087 Base Definitions volume of IEEE Std. 1003.1-200x, Table 5-1, Escape Sequences and Associated  
 2088 Actions, while the normal ERE behavior is to regard such a sequence as consisting of two  
 2089 characters. Allowing the *awk/lex* behavior in EREs would change the normal behavior in an  
 2090 unacceptable way; it is expected that *awk* and *lex* will decode escape sequences in EREs before  
 2091 passing them to *regcomp()* or comparable routines. Each utility describes the escape sequences it  
 2092 accepts as an exception to the rules in this section; the list is not the same, for historical reasons.

2093 As noted previously, the new syntax and rules have been added to accommodate other  
 2094 languages than English. The remainder of this section describes the rationale for these  
 2095 modifications.

2096 **A.9.3.6 BREs Matching Multiple Characters**

2097 The limit of nine back-references to subexpressions in the RE is based on the use of a single-digit  
 2098 identifier; increasing this to multiple digits would break historical applications. This does not  
 2099 imply that only nine subexpressions are allowed in REs. The following is a valid BRE with ten  
 2100 subexpressions:

2101 `\\(\\(\\(ab\\)*c\\)*d\\)\\(ef\\)*\\(gh\\){2}\\(ij\\)*\\(kl\\)*\\(mn\\)*\\(op\\)*\\(qr\\)*`



2102 The standard developers regarded the common historical behavior, which supported "\n\*", but  
 2103 not "\n\{min,max\}", "\(...\)\*", or "\(...\)\{min,max\}", as a non-intentional  
 2104 result of a specific implementation, and they supported both duplication and interval  
 2105 expressions following subexpressions and back-references.

2106 The changes to the processing of the back-reference expression remove an unspecified or  
 2107 ambiguous behavior in the Shell and Utilities volume of IEEE Std. 1003.1-200x, aligning it with  
 2108 the requirements specified for the *regcomp()* expression, and is the result of PASC Interpretation  
 2109 1003.2-92 #43 submitted for the ISO POSIX-2: 1993 standard.

#### 2110 A.9.3.7 BRE Precedence

2111 There is no additional rationale for this section.

#### 2112 A.9.3.8 BRE Expression Anchoring

2113 Often, the dollar sign is viewed as matching the ending <newline> in text files. This is not  
 2114 strictly true; the <newline> is typically eliminated from the strings to be matched, and the dollar  
 2115 sign matches the terminating null character.

2116 The ability of '^', '\$', and '\*' to be non-special in certain circumstances may be confusing to  
 2117 some programmers, but this situation was changed only in a minor way from historical practice  
 2118 to avoid breaking many historical scripts. Some consideration was given to making the use of  
 2119 the anchoring characters undefined if not escaped and not at the beginning or end of strings.  
 2120 This would cause a number of historical BREs, such as "2^10", "\$HOME", and "\$1.35", that  
 2121 relied on the characters being treated literally, to become invalid.

2122 However, one relatively uncommon case was changed to allow an extension used on some  
 2123 implementations. Historically, the BREs "^foo" and "\(^foo\)" did not match the same  
 2124 string, despite the general rule that subexpressions and entire BREs match the same strings. To  
 2125 increase consensus, IEEE Std. 1003.1-200x has allowed an extension on some systems to treat  
 2126 these two cases in the same way by declaring that anchoring *may* occur at the beginning or end  
 2127 of a subexpression. Therefore, portable BREs that require a literal circumflex at the beginning or  
 2128 a dollar sign at the end of a subexpression must escape them. Note that a BRE such as  
 2129 "a\(^bc\)" will either match "a^bc" or nothing on different systems under the rules.

2130 ERE anchoring has been different from BRE anchoring in all historical systems. An unescaped  
 2131 anchor character has never matched its literal counterpart outside a bracket expression. Some  
 2132 systems treated "foo\$bar" as a valid expression that never matched anything; others treated it  
 2133 as invalid. IEEE Std. 1003.1-200x mandates the former, valid unmatched behavior.

2134 Some systems have extended the BRE syntax to add alternation. For example, the subexpression  
 2135 "\(foo\$|bar\)" would match either "foo" at the end of the string or "bar" anywhere. The  
 2136 extension is triggered by the use of the undefined "\|" sequence. Because the BRE is undefined  
 2137 for portable scripts, the extending system is free to make other assumptions, such that the '\$'  
 2138 represents the end-of-line anchor in the middle of a subexpression. If it were not for the  
 2139 extension, the '\$' would match a literal dollar sign under the rules.

**2140 A.9.4 Extended Regular Expressions**

2141 As with BREs, the standard developers decided to make the interpretation of escaped ordinary  
2142 characters undefined.

2143 The right parenthesis is not listed as an ERE special character because it is only special in the  
2144 context of a preceding left parenthesis. If found without a preceding left parenthesis, the right  
2145 parenthesis has no special meaning.

2146 The *interval expression*, " $\{m,n\}$ ", has been added to EREs. Historically, the interval expression  
2147 has only been supported in some ERE implementations. The standard developers estimated that  
2148 the addition of interval expressions to EREs would not decrease consensus and would also make  
2149 BREs more of a subset of EREs than in many historical implementations.

2150 It was suggested that, in addition to interval expressions, back-references ( $\backslash n$ ) should also be  
2151 added to EREs. This was rejected by the standard developers as likely to decrease consensus.

2152 In historical implementations, multiple duplication symbols are usually interpreted from left to  
2153 right and treated as additive. As an example, " $a^*b$ " matches zero or more instances of 'a'  
2154 followed by a 'b'. In IEEE Std. 1003.1-200x, multiple duplication symbols are undefined; that is,  
2155 they cannot be relied upon for portable applications. One reason for this is to provide some  
2156 scope for future enhancements.

2157 The precedence of operations differs between EREs and those in *lex*; in *lex*, for historical reasons,  
2158 interval expressions have a lower precedence than concatenation.

**2159 A.9.4.1 EREs Matching a Single Character or Collating Element**

2160 There is no additional rationale for this section.

**2161 A.9.4.2 ERE Ordinary Characters**

2162 There is no additional rationale for this section.

**2163 A.9.4.3 ERE Special Characters**

2164 There is no additional rationale for this section.

**2165 A.9.4.4 Periods in EREs**

2166 There is no additional rationale for this section.

**2167 A.9.4.5 ERE Bracket Expression**

2168 There is no additional rationale for this section.

**2169 A.9.4.6 EREs Matching Multiple Characters**

2170 There is no additional rationale for this section.

**2171 A.9.4.7 ERE Alternation**

2172 There is no additional rationale for this section.

2173 A.9.4.8 *ERE Precedence*

2174 There is no additional rationale for this section.

2175 A.9.4.9 *ERE Expression Anchoring*

2176 There is no additional rationale for this section.

2177 **A.9.5 Regular Expression Grammar**

2178 The grammars are intended to represent the range of acceptable syntaxes available to portable  
 2179 applications. There are instances in the text where undefined constructs are described; as  
 2180 explained previously, these allow implementation extensions. There is no intended requirement  
 2181 that an implementation extension must somehow fit into the grammars shown here.

2182 The BRE grammar does not permit L\_ANCHOR or R\_ANCHOR inside "\(" and "\)" (which  
 2183 implies that '^' and '\$' are ordinary characters). This reflects the semantic limits on the  
 2184 application, as noted in the Base Definitions volume of IEEE Std. 1003.1-200x, Section 9.3.8, BRE  
 2185 Expression Anchoring. Implementations are permitted to extend the language to interpret '^'  
 2186 and '\$' as anchors in these locations, and as such, portable applications cannot use unescaped  
 2187 '^' and '\$' in positions inside "\(" and "\)" that might be interpreted as anchors.

2188 The ERE grammar does not permit several constructs that the Base Definitions volume of  
 2189 IEEE Std. 1003.1-200x, Section 9.4.2, ERE Ordinary Characters and the Base Definitions volume  
 2190 of IEEE Std. 1003.1-200x, Section 9.4.3, ERE Special Characters specify as having undefined  
 2191 results:

- 2192 • ORD\_CHAR preceded by '\'
- 2193 • *ERE\_dupl\_symbol*(s) appearing first in an ERE, or immediately following '|', '^', or '('
- 2194 • '{' not part of a valid *ERE\_dupl\_symbol*
- 2195 • '|' appearing first or last in an ERE, or immediately following '|' or '(', or immediately  
 2196 preceding ')'

2197 Implementations are permitted to extend the language to allow these. Portable applications  
 2198 cannot use such constructs.

2199 A.9.5.1 *BRE/ERE Grammar Lexical Conventions*

2200 There is no additional rationale for this section.

2201 A.9.5.2 *RE and Bracket Expression Grammar*

2202 The removal of the *Back\_open\_paren Back\_close\_paren* option from the *nondupl\_RE* specification is  
 2203 the result of PASC Interpretation 1003.2-92 #43 submitted for the ISO POSIX-2: 1993 standard.  
 2204 Although the grammar required support for null subexpressions, this section does not describe  
 2205 the meaning of, and historical practice did not support, this construct.

2206 A.9.5.3 *ERE Grammar*

2207 There is no additional rationale for this section.

## 2208 **A.10 Directory Structure and Devices**

### 2209 **A.10.1 Directory Structure and Files**

2210 A description of the historical **/usr/tmp** was omitted, removing any concept of differences in  
2211 emphasis between the **/** and **/usr** directories. The descriptions of **/bin**, **/usr/bin**, **/lib**, and **/usr/lib**  
2212 were omitted because they are not useful for applications. In an early draft, a distinction was  
2213 made between system and application directory usage, but this was not found to be useful.

2214 The directories **/** and **/dev** are included because the notion of a hierarchical directory structure is  
2215 key to other information presented elsewhere in IEEE Std. 1003.1-200x. In early drafts, it was  
2216 argued that special devices and temporary files could conceivably be handled without a  
2217 directory structure on some implementations. For example, the system could treat the characters  
2218 `" /tmp "` as a special token that would store files using some non-POSIX file system structure.  
2219 This notion was rejected by the standard developers, who required that all the files in this  
2220 section be implemented via POSIX file systems.

2221 The **/tmp** directory is retained in IEEE Std. 1003.1-200x to accommodate historical applications  
2222 that assume its availability. Implementations are encouraged to provide suitable directory  
2223 names in the environment variable *TMPDIR* and applications are encouraged to use the contents  
2224 of *TMPDIR* for creating temporary files.

2225 The standard files **/dev/null** and **/dev/tty** are required to be both readable and writable to allow  
2226 applications to have the intended historical access to these files.

2227 The standard file **/dev/console** has been added for alignment with the Single UNIX Specification.

### 2228 **A.10.2 Output Devices and Terminal Types**

2229 There is no additional rationale for this section.

## 2230 A.11 General Terminal Interface

2231 If the implementation does not support this interface on any device types, it should behave as if  
2232 it were being used on a device that is not a terminal device (in most cases *errno* will be set to  
2233 [ENOTTY] on return from functions defined by this interface). This is based on the fact that  
2234 many applications are written to run both interactively and in some non-interactive mode, and  
2235 they adapt themselves at runtime. Requiring that they all be modified to test an environment  
2236 variable to determine whether they should try to adapt is unnecessary. On a system that  
2237 provides no general terminal interface, providing all the entry points as stubs that return  
2238 [ENOTTY] (or an equivalent, as appropriate) has the same effect and requires no changes to the  
2239 application.

2240 Although the needs of both interface implementors and application developers were addressed  
2241 throughout IEEE Std. 1003.1-200x, this section pays more attention to the needs of the latter. This  
2242 is because, while many aspects of the programming interface can be hidden from the user by the  
2243 application developer, the terminal interface is usually a large part of the user interface.  
2244 Although to some extent the application developer can build missing features or work around  
2245 inappropriate ones, the difficulties of doing that are greater in the terminal interface than  
2246 elsewhere. For example, efficiency prohibits the average program from interpreting every  
2247 character passing through it in order to simulate character erase, line kill, and so on. These  
2248 functions should usually be done by the operating system, possibly at the interrupt level.

2249 The *tc\**() functions were introduced as a way of avoiding the problems inherent in the  
2250 traditional *ioctl*() function and in variants of it that were proposed. For example, *tcsetattr*() is  
2251 specified in place of the use of the TCSETA *ioctl*() command function. This allows specification  
2252 of all the arguments in a manner consistent with the ISO C standard unlike the varying third  
2253 argument of *ioctl*(), which is sometimes a pointer (to any of many different types) and  
2254 sometimes an **int**.

2255 The advantages of this new method include:

- 2256 • It allows strict type checking.
- 2257 • The direction of transfer of control data is explicit.
- 2258 • Portable capabilities are clearly identified.
- 2259 • The need for a general interface routine is avoided.
- 2260 • Size of the argument is well-defined (there is only one type).

2261 The disadvantages include:

- 2262 • No historical implementation uses the new method.
- 2263 • There are many small routines instead of one general-purpose one.
- 2264 • The historical parallel with *fcntl*() is broken.

2265 The issue of modem control was excluded from IEEE Std. 1003.1-200x on the grounds that:

- 2266 • It was concerned with setting and control of hardware timers.
- 2267 • The appropriate timers and settings vary widely internationally.
- 2268 • Feedback from European computer manufacturers indicated that this facility was not  
2269 consistent with European needs and that specification of such a facility was not a  
2270 requirement for portability.

2271 **A.11.1 Interface Characteristics**2272 *A.11.1.1 Opening a Terminal Device File*

2273 There is no additional rationale provided for this section.

2274 *A.11.1.2 Process Groups*

2275 There is a potential race when the members of the foreground process group on a terminal leave  
2276 that process group, either by exit or by changing process groups. After the last process exits the  
2277 process group, but before the foreground process group ID of the terminal is changed (usually  
2278 by a job-control shell), it would be possible for a new process to be created with its process ID  
2279 equal to the terminal's foreground process group ID. That process might then become the  
2280 process group leader and accidentally be placed into the foreground on a terminal that was not  
2281 necessarily its controlling terminal. As a result of this problem, the controlling terminal is  
2282 defined to not have a foreground process group during this time.

2283 The cases where a controlling terminal has no foreground process group occur when all  
2284 processes in the foreground process group either terminate and are waited for or join other  
2285 process groups via *setpgid()* or *setsid()*. If the process group leader terminates, this is the first  
2286 case described; if it leaves the process group via *setpgid()*, this is the second case described (a  
2287 process group leader cannot successfully call *setsid()*). When one of those cases causes a  
2288 controlling terminal to have no foreground process group, it has two visible effects on  
2289 applications. The first is the value returned by *tcgetpgrp()*. The second (which occurs only in the  
2290 case where the process group leader terminates) is the sending of signals in response to special  
2291 input characters. The intent of IEEE Std. 1003.1-200x is that no process group be wrongly  
2292 identified as the foreground process group by *tcgetpgrp()* or unintentionally receive signals  
2293 because of placement into the foreground.

2294 In 4.3 BSD, the old process group ID continues to be used to identify the foreground process  
2295 group and is returned by the function equivalent to *tcgetpgrp()*. In that implementation it is  
2296 possible for a newly created process to be assigned the same value as a process ID and then form  
2297 a new process group with the same value as a process group ID. The result is that the new  
2298 process group would receive signals from this terminal for no apparent reason, and  
2299 IEEE Std. 1003.1-200x precludes this by forbidding a process group from entering the foreground  
2300 in this way. It would be more direct to place part of the requirement made by the last sentence  
2301 under *fork()*, but there is no convenient way for that section to refer to the value that *tcgetpgrp()*  
2302 returns, since in this case there is no process group and thus no process group ID.

2303 One possibility for a conforming implementation is to behave similarly to 4.3 BSD, but to  
2304 prevent this reuse of the ID, probably in the implementation of *fork()*, as long as it is in use by  
2305 the terminal.

2306 Another possibility is to recognize when the last process stops using the terminal's foreground  
2307 process group ID, which is when the process group lifetime ends, and to change the terminal's  
2308 foreground process group ID to a reserved value that is never used as a process ID or process  
2309 group ID. (See the definition of *process group lifetime* in the definitions section.) The process ID  
2310 can then be reserved until the terminal has another foreground process group.

2311 The 4.3 BSD implementation permits the leader (and only member) of the foreground process  
2312 group to leave the process group by calling the equivalent of *setpgid()* and to later return,  
2313 expecting to return to the foreground. There are no known application needs for this behavior,  
2314 and IEEE Std. 1003.1-200x neither requires nor forbids it (except that it is forbidden for session  
2315 leaders) by leaving it unspecified.

2316 A.11.1.3 *The Controlling Terminal*

2317 IEEE Std. 1003.1-200x does not specify a mechanism by which to allocate a controlling terminal.  
2318 This is normally done by a system utility (such as *getty*) and is considered an administrative  
2319 feature outside the scope of IEEE Std. 1003.1-200x.

2320 Historical implementations allocate controlling terminals on certain *open()* calls. Since *open()* is  
2321 part of POSIX.1, its behavior had to be dealt with. The traditional behavior is not required  
2322 because it is not very straightforward or flexible for either implementations or applications.  
2323 However, because of its prevalence, it was not practical to disallow this behavior either. Thus, a  
2324 mechanism was standardized to ensure portable, predictable behavior in *open()*.

2325 Some historical implementations deallocate a controlling terminal on the last system-wide close.  
2326 This behavior is neither required nor prohibited. Even on implementations that do provide this  
2327 behavior, applications generally cannot depend on it due to its system-wide nature.

2328 A.11.1.4 *Terminal Access Control*

2329 The access controls described in this section apply only to a process that is accessing its  
2330 controlling terminal. A process accessing a terminal that is not its controlling terminal is  
2331 effectively treated the same as a member of the foreground process group. While this may seem  
2332 unintuitive, note that these controls are for the purpose of job control, not security, and job  
2333 control relates only to a process' controlling terminal. Normal file access permissions handle  
2334 security.

2335 If the process calling *read()* or *write()* is in a background process group that is orphaned, it is not  
2336 desirable to stop the process group, as it is no longer under the control of a job control shell that  
2337 could put it into foreground again. Accordingly, calls to *read()* or *write()* functions by such  
2338 processes receive an immediate error return. This is different than in 4.2 BSD, which kills  
2339 orphaned processes that receive terminal stop signals.

2340 The foreground/background/orphaned process group check performed by the terminal driver  
2341 must be repeatedly performed until the calling process moves into the foreground or until the  
2342 process group of the calling process becomes orphaned. That is, when the terminal driver  
2343 determines that the calling process is in the background and should receive a job control signal,  
2344 it sends the appropriate signal (SIGTTIN or SIGTTOU) to every process in the process group of  
2345 the calling process and then it allows the calling process to immediately receive the signal. The  
2346 latter is typically performed by blocking the process so that the signal is immediately noticed.  
2347 Note, however, that after the process finishes receiving the signal and control is returned to the  
2348 driver, the terminal driver must reexecute the foreground/background/orphaned process group  
2349 check. The process may still be in the background, either because it was continued in the  
2350 background by a job-control shell, or because it caught the signal and did nothing.

2351 The terminal driver repeatedly performs the foreground/background/orphaned process group  
2352 checks whenever a process is about to access the terminal. In the case of *write()* or the control  
2353 *tc\*()* functions, the check is performed at the entry of the function. In the case of *read()*, the check  
2354 is performed not only at the entry of the function, but also after blocking the process to wait for  
2355 input characters (if necessary). That is, once the driver has determined that the process calling  
2356 the *read()* function is in the foreground, it attempts to retrieve characters from the input queue. If  
2357 the queue is empty, it blocks the process waiting for characters. When characters are available  
2358 and control is returned to the driver, the terminal driver must return to the repeated  
2359 foreground/background/orphaned process group check again. The process may have moved  
2360 from the foreground to the background while it was blocked waiting for input characters.

2361 A.11.1.5 *Input Processing and Reading Data*

2362 There is no additional rationale provided for this section.

2363 A.11.1.6 *Canonical Mode Input Processing*

2364 The term *character* is intended here. ERASE should erase the last character, not the last byte. In  
2365 the case of multi-byte characters, these two may be different.

2366 4.3 BSD has a WERASE character that erases the last “word” typed (but not any preceding  
2367 <blank>s or <tab>s). A word is defined as a sequence of non-<blank> characters, with <tab>s  
2368 counted as <blank>s. Like ERASE, WERASE does not erase beyond the beginning of the line.  
2369 This WERASE feature has not been specified in POSIX.1 because it is difficult to define in the  
2370 international environment. It is only useful for languages where words are delimited by  
2371 <blank>s. In some ideographic languages, such as Japanese and Chinese, words are not  
2372 delimited at all. The WERASE character should presumably take one back to the beginning of a  
2373 sentence in those cases; practically, this means it would not get much use for those languages.

2374 It should be noted that there is a possible inherent deadlock if the application and  
2375 implementation conflict on the value of MAX\_CANON. With ICANON set (if IXOFF is  
2376 enabled) and more than MAX\_CANON characters transmitted without a <linefeed>,  
2377 transmission will be stopped, the <linefeed> (or <carriage-return> when ICRLF is set) will never  
2378 arrive, and the *read()* will never be satisfied.

2379 An application should not set IXOFF if it is using canonical mode unless it knows that (even in  
2380 the face of a transmission error) the conditions described previously cannot be met or unless it is  
2381 prepared to deal with the possible deadlock in some other way, such as timeouts.

2382 It should also be noted that this can be made to happen in non-canonical mode if the trigger  
2383 value for sending IXOFF is less than VMIN and VTIME is zero.

2384 A.11.1.7 *Non-Canonical Mode Input Processing*

2385 Some points to note about MIN and TIME:

- 2386 1. The interactions of MIN and TIME are not symmetric. For example, when MIN>0 and  
2387 TIME=0, TIME has no effect. However, in the opposite case where MIN=0 and TIME>0,  
2388 both MIN and TIME play a role in that MIN is satisfied with the receipt of a single  
2389 character.
- 2390 2. Also note that in case A (MIN>0, TIME>0), TIME represents an inter-character timer, while  
2391 in case C (MIN=0, TIME>0), TIME represents a read timer.

2392 These two points highlight the dual purpose of the MIN/TIME feature. Cases A and B, where  
2393 MIN>0, exist to handle burst-mode activity (for example, file transfer programs) where a  
2394 program would like to process at least MIN characters at a time. In case A, the inter-character  
2395 timer is activated by a user as a safety measure; in case B, it is turned off.

2396 Cases C and D exist to handle single-character timed transfers. These cases are readily adaptable  
2397 to screen-based applications that need to know if a character is present in the input queue before  
2398 refreshing the screen. In case C, the read is timed; in case D, it is not.

2399 Another important note is that MIN is always just a minimum. It does not denote a record  
2400 length. That is, if a program does a read of 20 bytes, MIN is 10, and 25 characters are present, 20  
2401 characters shall be returned to the user. In the special case of MIN=0, this still applies: if more  
2402 than one character is available, they all will be returned immediately.



2403 **A.11.1.8 Writing Data and Output Processing**

2404 There is no additional rationale for this section.

2405 **A.11.1.9 Special Characters**

2406 There is no additional rationale for this section.

2407 **A.11.1.10 Modem Disconnect**

2408 There is no additional rationale for this section.

2409 **A.11.1.11 Closing a Terminal Device File**

2410 IEEE Std. 1003.1-200x does not specify that a *close()* on a terminal device file include the  
2411 equivalent of a call to *tcfow(fd,TCOON)*.

2412 An implementation that discards output at the time *close()* is called after reporting the return  
2413 value to the *write()* call that data was written does not conform with IEEE Std. 1003.1-200x. An  
2414 application has functions such as *tcdrain()*, *tcfush()*, and *tcfow()* available to obtain the detailed  
2415 behavior it requires with respect to flushing of output.

2416 At the time of the last close on a terminal device, an application relinquishes any ability to exert  
2417 flow control via *tcfow()*.

2418 **A.11.2 Parameters that Can be Set**2419 **A.11.2.1 The termios Structure**

2420 This structure is part of an interface that, in general, retains the historic grouping of flags.  
2421 Although a more optimal structure for implementations may be possible, the degree of change  
2422 to applications would be significantly larger.

2423 **A.11.2.2 Input Modes**

2424 Some historical implementations treated a long break as multiple events, as many as one per  
2425 character time. The wording in POSIX.1 explicitly prohibits this.

2426 Although the ISTRIP flag is normally superfluous with today's terminal hardware and software,  
2427 it is historically supported. Therefore, applications may be using ISTRIP, and there is no  
2428 technical problem with supporting this flag. Also, applications may wish to receive only 7-bit  
2429 input bytes and may not be connected directly to the hardware terminal device (for example,  
2430 when a connection traverses a network).

2431 Also, there is no requirement in general that the terminal device ensures that high-order bits  
2432 beyond the specified character size are cleared. ISTRIP provides this function for 7-bit  
2433 characters, which are common.

2434 In dealing with multi-byte characters, the consequences of a parity error in such a character, or in  
2435 an escape sequence affecting the current character set, are beyond the scope of POSIX.1 and are  
2436 best dealt with by the application processing the multi-byte characters.

2437 *A.11.2.3 Output Modes*

2438 POSIX.1 does not describe postprocessing of output to a terminal or detailed control of that from  
2439 a portable application. (That is, translation of <newline> to <carriage-return> followed by  
2440 <linefeed> or <tab> processing.) There is nothing that a portable application should do to its  
2441 output for a terminal because that would require knowledge of the operation of the terminal. It  
2442 is the responsibility of the operating system to provide postprocessing appropriate to the output  
2443 device, whether it is a terminal or some other type of device.

2444 Extensions to POSIX.1 to control the type of postprocessing already exist and are expected to  
2445 continue into the future. The control of these features is primarily to adjust the interface between  
2446 the system and the terminal device so the output appears on the display correctly. This should  
2447 be set up before use by any application.

2448 In general, both the input and output modes should not be set absolutely, but rather modified  
2449 from the inherited state.

2450 *A.11.2.4 Control Modes*

2451 This section could be misread that the symbol “CSIZE” is a title in the **termios** *c\_flag* field .  
2452 Although it does serve that function, it is also a required symbol, as a literal reading of POSIX.1  
2453 (and the caveats about typography) would indicate.

2454 *A.11.2.5 Local Modes*

2455 Non-canonical mode is provided to allow fast bursts of input to be read efficiently while still  
2456 allowing single-character input.

2457 The ECHONL function historically has been in many implementations. Since there seems to be  
2458 no technical problem with supporting ECHONL, it is included in POSIX.1 to increase consensus.

2459 The alternate behavior possible when ECHOK or ECHOE are specified with ICANON is  
2460 permitted as a compromise depending on what the actual terminal hardware can do. Erasing  
2461 characters and lines is preferred, but is not always possible.

2462 *A.11.2.6 Special Control Characters*

2463 Permitting VMIN and VTIME to overlap with VEOF and VEOL was a compromise for historical  
2464 implementations. Only when backwards-compatibility of object code is a serious concern to an  
2465 implementor should an implementation continue this practice. Correct applications that work  
2466 with the overlap (at the source level) should also work if it is not present, but not the reverse.

2467 **A.12 Utility Conventions**2468 **A.12.1 Utility Argument Syntax**

2469 The standard developers considered that recent trends toward diluting the SYNOPSIS sections  
2470 of historical reference pages to the equivalent of:

2471 `command [options][operands]`

2472 were a disservice to the reader. Therefore, considerable effort was placed into rigorous  
2473 definitions of all the command line arguments and their interrelationships. The relationships  
2474 depicted in the synopses are normative parts of IEEE Std. 1003.1-200x; this information is  
2475 sometimes repeated in textual form, but that is only for clarity within context.

2476 The use of “undefined” for conflicting argument usage and for repeated usage of the same  
2477 option is meant to prevent portable applications from using conflicting arguments or repeated  
2478 options unless specifically allowed (as is the case with *ls*, which allows simultaneous, repeated  
2479 use of the `-C`, `-l`, and `-1` options). Many historical implementations will tolerate this usage,  
2480 choosing either the first or the last applicable argument. This tolerance can continue, but  
2481 portable applications cannot rely upon it. (Other implementations may choose to print usage  
2482 messages instead.)

2483 The use of “undefined” for conflicting argument usage also allows an implementation to make  
2484 reasonable extensions to utilities where the implementor considers mutually-exclusive options  
2485 according to IEEE Std. 1003.1-200x to have a sensible meaning and result.

2486 IEEE Std. 1003.1-200x does not define the result of a command when an option-argument or  
2487 operand is not followed by ellipses and the application specifies more than one of that option-  
2488 argument or operand. This allows an implementation to define valid (although non-standard)  
2489 behavior for the utility when more than one such option or operand is specified.

2490 Allowing <blank> characters after an option (that is, placing an option and its option-argument  
2491 into separate argument strings) when IEEE Std. 1003.1-200x does not require it encourages  
2492 portability of users, while still preserving backwards-compatibility of scripts. Inserting <blank>  
2493 characters between the option and the option-argument is preferred; however, historical usage  
2494 has not been consistent in this area; therefore, <blank>s are required to be handled by all  
2495 implementations, but implementations are also allowed to handle the historical syntax. Another  
2496 justification for selecting the multiple-argument method was that the single-argument case is  
2497 inherently ambiguous when the option-argument can legitimately be a null string.

2498 IEEE Std. 1003.1-200x explicitly states that digits are permitted as operands and option-  
2499 arguments. The lower and upper bounds for the values of the numbers used for operands and  
2500 option-arguments were derived from the ISO C standard values for `{LONG_MIN}` and  
2501 `{LONG_MAX}`. The requirement on the standard utilities is that numbers in the specified range  
2502 do not cause a syntax error, although the specification of a number need not be semantically  
2503 correct for a particular operand or option-argument of a utility. For example, the specification of:

2504 `dd obs=3000000000`

2505 would yield undefined behavior for the application and would be a syntax error because the  
2506 number 3 000 000 000 is outside of the range `-2 147 483 647` to `+2 147 483 647`. On the other hand:

2507 `dd obs=2000000000`

2508 may cause some error, such as “blocksize too large”, rather than a syntax error.

2509 **A.12.2 Utility Syntax Guidelines**

2510 This section is based on the rules listed in the SVID. It was included for two reasons:

- 2511 1. The individual utility descriptions in the Shell and Utilities volume of  
2512 IEEE Std. 1003.1-200x, Chapter 4, Utilities needed a set of common (although not universal)  
2513 actions on which they could anchor their descriptions of option and operand syntax. Most  
2514 of the standard utilities actually do use these guidelines, and many of their historical  
2515 implementations use the *getopt()* function for their parsing. Therefore, it was simpler to  
2516 cite the rules and merely identify exceptions.
- 2517 2. Writers of portable applications need suggested guidelines if the POSIX community is to  
2518 avoid the chaos of historical UNIX system command syntax.

2519 It is recommended that all *future* utilities and applications use these guidelines to enhance “user  
2520 portability”. The fact that some historical utilities could not be changed (to avoid breaking  
2521 historical applications) should not deter this future goal.

2522 The voluntary nature of the guidelines is highlighted by repeated uses of the word *should*  
2523 throughout. This usage should not be misinterpreted to imply that utilities that claim  
2524 conformance in their OPTIONS sections do not always conform.

2525 Guidelines 1 and 2 are offered as guidance for locales using Latin alphabets. No  
2526 recommendations are made by IEEE Std. 1003.1-200x concerning utility naming in other locales.

2527 In the Shell and Utilities volume of IEEE Std. 1003.1-200x, Section 2.9.1, Simple Commands, it is  
2528 further stated that a command used in the Shell Command Language cannot be named with a  
2529 trailing colon.

2530 Guideline 3 was changed to allow alphanumeric characters (letters and digits) from the character  
2531 set to allow compatibility with historical usage. Historical practice allows the use of digits  
2532 wherever practical, and there are no portability issues that would prohibit the use of digits. In  
2533 fact, from an internationalization viewpoint, digits (being non-language-dependent) are  
2534 preferable over letters (a *-2* is intuitively self-explanatory to any user, while in the *-f filename* the  
2535 letter ‘f’ is a mnemonic aid only to speakers of Latin-based languages where “file name”  
2536 happens to translate to a word that begins with ‘f’). Since guideline 3 still retains the word  
2537 “single”, multi-digit options are not allowed. Instances of historical utilities that used them have  
2538 been marked obsolescent, with the numbers being changed from option names to option-  
2539 arguments.

2540 It was difficult to achieve a satisfactory solution to the problem of name space in option  
2541 characters. When the standard developers desired to extend the historical *cc* utility to accept  
2542 ISO C standard programs, they found that all of the portable alphabet was already in use by  
2543 various vendors. Thus, they had to devise a new name, *c89*, rather than something like *cc -X*.  
2544 There were suggestions that implementors be restricted to providing extensions through various  
2545 means (such as using a plus sign as the option delimiter or using option characters outside the  
2546 alphanumeric set) that would reserve all of the remaining alphanumeric characters for future  
2547 POSIX standards. These approaches were resisted because they lacked the historical style of  
2548 UNIX systems. Furthermore, if a vendor-provided option should become commonly used in the  
2549 industry, it would be a candidate for standardization. It would be desirable to standardize such a  
2550 feature using historical practice for the syntax (the semantics can be standardized with any  
2551 syntax). This would not be possible if the syntax was one reserved for the vendor. However,  
2552 since the standardization process may lead to minor changes in the semantics, it may prove to be  
2553 better for a vendor to use a syntax that will not be affected by standardization.

2554 Guideline 8 includes the concept of comma-separated lists in a single argument. It is up to the  
2555 utility to parse such a list itself because *getopt()* just returns the single string. This situation was

2556 retained so that certain historical utilities would not violate the guidelines. Applications  
2557 preparing for international use should be aware of an occasional problem with comma-  
2558 separated lists: in some locales, the comma is used as the radix character. Thus, if an application  
2559 is preparing operands for a utility that expects a comma-separated lists, it should avoid  
2560 generating non-integer values through one of the means that is influenced by setting the  
2561 *LC\_NUMERIC* variable (such as *awk*, *bc*, *printf*, or *printf()*).

2562 Applications calling any utility with a first operand starting with '-' should usually specify --,  
2563 as indicated by Guideline 10, to mark the end of the options. This is true even if the SYNOPSIS in  
2564 the Shell and Utilities volume of IEEE Std. 1003.1-200x does not specify any options;  
2565 implementations may provide options as extensions to the Shell and Utilities volume of  
2566 IEEE Std. 1003.1-200x. The standard utilities that do not support Guideline 10 indicate that fact  
2567 in the OPTIONS section of the utility description.

2568 Guideline 11 was modified to clarify that the order of different options should not matter  
2569 relative to one another. However, the order of repeated options that also have option-arguments  
2570 may be significant; therefore, such options are required to be interpreted in the order that they  
2571 are specified. The *make* utility is an instance of a historical utility that uses repeated options  
2572 in which the order is significant. Multiple files are specified by giving multiple instances of the -f  
2573 option; for example:

```
2574     make -f common_header -f specific_rules target
```

2575 Guideline 13 does not imply that all of the standard utilities automatically accept the operand  
2576 '-' to mean standard input or output, nor does it specify the actions of the utility upon  
2577 encountering multiple '-' operands. It simply says that, by default, '-' operands are not used  
2578 for other purposes in the file reading or writing (but not when using *stat*, *unlink*, *touch*, and so on)  
2579 utilities. All information concerning actual treatment of the '-' operand is found in the  
2580 individual utility sections.

2581 An area of concern was that as implementations mature, implementation-defined utilities and  
2582 implementation-defined utility options will result. The idea was expressed that there needed to  
2583 be a standard way, say an environment variable or some such mechanism, to identify  
2584 implementation-defined utilities separately from standard utilities that may have the same  
2585 name. It was decided that there already exist several ways of dealing with this situation and that  
2586 it is outside of the POSIX.2 scope to attempt to standardize in the area of non-standard items. A  
2587 method that exists on some historical implementations is the use of the so-called */local/bin* or  
2588 */usr/local/bin* directory to separate local or additional copies or versions of utilities. Another  
2589 method that is also used is to isolate utilities into completely separate domains. Still another  
2590 method to ensure that the desired utility is being used is to request the utility by its full path  
2591 name. There are many approaches to this situation; the examples given above serve to illustrate  
2592 that there is more than one.

**2593 A.13 Headers****2594 A.13.1 Format of Entries**

2595 Each header reference page has a common layout of sections describing the interface. This layout  
2596 is similar to the manual page or “man” page format shipped with most UNIX systems, and each  
2597 header has sections describing the SYNOPSIS and DESCRIPTION. These are the two sections  
2598 that relate to conformance.

2599 Additional sections are informative, and add considerable information for the application  
2600 developer. APPLICATION USAGE sections provide additional caveats, issues, and  
2601 recommendations to the developer. RATIONALE sections give additional information on the  
2602 decisions made in defining the interface.

2603 FUTURE DIRECTIONS sections act as pointers to related work that may impact the interface in  
2604 the future, and often cautions the developer to architect the code to account for a change in this  
2605 area. Note that a future directions statement should not be taken as a commitment to adopt a  
2606 feature or interface in the future.

2607 The CHANGE HISTORY section describes when the interface was introduced, and how it has  
2608 changed.

2609 Option labels and margin markings in the page can be useful in guiding the application  
2610 developer.

2611 / *Rationale (Informative)*

2612 **Part B:**

2613 **System Interfaces**

2614 *The Open Group*





# Rationale for System Interfaces

2615

## 2616 **B.1 Introduction**

### 2617 **B.1.1 Scope**

2618 Refer to Section A.1.1 (on page 3311).

### 2619 **B.1.2 Conformance**

2620 Refer to Section A.2 (on page 3317).

### 2621 **B.1.3 Normative References**

2622 There is no additional rationale for this section.

### 2623 **B.1.4 Changes from Issue 4**

2624 The change history is provided as an informative section, to track changes from previous issues  
2625 of IEEE Std. 1003.1-200x that comprised earlier versions of the Single UNIX Specification.

#### 2626 *B.1.4.1 Changes from Issue 4 to Issue 4, Version 2*

2627 There is no additional rationale for this section.

#### 2628 *B.1.4.2 Changes from Issue 4, Version 2 to Issue 5*

2629 There is no additional rationale for this section.

#### 2630 *B.1.4.3 Changes from Issue 5 to Issue 6 (IEEE Std. 1003.1-200x)*

2631 There is no additional rationale for this section.

### 2632 **B.1.5 New Features**

2633 There is no additional rationale for this section.

#### 2634 *B.1.5.1 New Features in Issue 4, Version 2*

2635 There is no additional rationale for this section.

#### 2636 *B.1.5.2 New Features in Issue 5*

2637 There is no additional rationale for this section.

#### 2638 *B.1.5.3 New Features in Issue 6*

2639 There is no additional rationale for this section.

**2640 B.1.6 Terminology**

2641 Refer to Section A.1.4 (on page 3313).

**2642 B.1.7 Definitions**

2643 Refer to Section A.3 (on page 3321).

**2644 B.1.8 Relationship to Other Formal Standards**

2645 There is no additional rationale for this section.

**2646 B.1.9 Portability**

2647 Refer to Section A.1.5 (on page 3315).

**2648 B.1.9.1 Codes**

2649 Refer to Section A.1.5.1 (on page 3315).

**2650 B.1.10 Format of Entries**

2651 Each system interface reference page has a common layout of sections describing the interface.  
2652 This layout is similar to the manual page or “man” page format shipped with most UNIX  
2653 systems, and each header has sections describing the SYNOPSIS, DESCRIPTION, RETURN  
2654 VALUE, and ERRORS. These are the four sections that relate to conformance.

2655 Additional sections are informative, and add considerable information for the application  
2656 developer. EXAMPLES sections provide example usage. APPLICATION USAGE sections  
2657 provide additional caveats, issues, and recommendations to the developer. RATIONALE  
2658 sections give additional information on the decisions made in defining the interface.

2659 FUTURE DIRECTIONS sections act as pointers to related work that may impact the interface in  
2660 the future, and often cautions the developer to architect the code to account for a change in this  
2661 area. Note that a future directions statement should not be taken as a commitment to adopt a  
2662 feature or interface in the future.

2663 The CHANGE HISTORY section describes when the interface was introduced, and how it has  
2664 changed.

2665 Option labels and margin markings in the page can be useful in guiding the application  
2666 developer.

2667 **B.2 General Information**2668 **B.2.1 Use and Implementation of Functions**

2669 The information concerning the use of functions was adapted from a description in the ISO C  
 2670 standard. Here is an example of how an application program can protect itself from library  
 2671 functions that may or may not be macros, rather than true functions:

2672 The *atoi()* function may be used in any of several ways:

- 2673 • By use of its associated header (possibly generating a macro expansion):

```
2674     #include <stdlib.h>
2675     /* ... */
2676     i = atoi(str);
```

- 2677 • By use of its associated header (assuredly generating a true function call):

```
2678     #include <stdlib.h>
2679     #undef atoi
2680     /* ... */
2681     i = atoi(str);
```

2682 or:

```
2683     #include <stdlib.h>
2684     /* ... */
2685     i = (atoi) (str);
```

- 2686 • By explicit declaration:

```
2687     extern int atoi (const char *);
2688     /* ... */
2689     i = atoi(str);
```

- 2690 • By implicit declaration:

```
2691     /* ... */
2692     i = atoi(str);
```

2693 (Assuming no function prototype is in scope. This is not allowed by the ISO C standard for  
 2694 functions with variable arguments; furthermore, parameter type conversion “widening” is  
 2695 subject to different rules in this case.)

2696 Note that the ISO C standard reserves names starting with ‘\_’ for the compiler. Therefore, the  
 2697 compiler could, for example, implement an intrinsic, built-in function *\_asm\_builtin\_atoi()*, which  
 2698 it recognized and expanded into inline assembly code. Then, in *<stdlib.h>*, there could be the  
 2699 following:

```
2700     #define atoi(X) _asm_builtin_atoi(X)
```

2701 The user’s “normal” call to *atoi()* would then be expanded inline, but the implementor would  
 2702 also be required to provide a callable function named *atoi()* for use when the application  
 2703 requires it; for example, if its address is to be stored in a function pointer variable.

## 2704 **B.2.2 The Compilation Environment**

### 2705 *B.2.2.1 POSIX.1 Symbols*

2706 This and the following section address the issue of “name space pollution”. The ISO C standard  
2707 requires that the name space beyond what it reserves not be altered except by explicit action of  
2708 the application writer. This section defines the actions to add the POSIX.1 symbols for those  
2709 headers where both the ISO C standard and POSIX.1 need to define symbols, and also where the  
2710 XSI Extension extends the base standard.

2711 When headers are used to provide symbols, there is a potential for introducing symbols that the  
2712 application writer cannot predict. Ideally, each header should only contain one set of symbols,  
2713 but this is not practical for historical reasons. Thus, the concept of feature test macros is  
2714 included. Two feature test macros are explicitly defined by IEEE Std. 1003.1-200x; it is expected  
2715 that future revisions may add to this.

2716 It is further intended that these feature test macros apply only to the headers specified by  
2717 IEEE Std. 1003.1-200x. Implementations are expressly permitted to make visible symbols not  
2718 specified by IEEE Std. 1003.1-200x, within both POSIX.1 and other headers, under the control of  
2719 feature test macros that are not defined by IEEE Std. 1003.1-200x.

### 2720 **The `_POSIX_C_SOURCE` Feature Test Macro**

2721 Since `_POSIX_SOURCE` specified by the POSIX.1-1990 standard did not have a value associated  
2722 with it, the `_POSIX_C_SOURCE` macro replaces it, allowing an application to inform the system  
2723 of the revision of the standard to which it conforms. This symbol will allow implementations to  
2724 support various revisions of IEEE Std. 1003.1-200x simultaneously. For instance, when either  
2725 `_POSIX_SOURCE` is defined or `_POSIX_C_SOURCE` is defined as 1, the system should make  
2726 visible the same name space as permitted and required by the POSIX.1-1990 standard. When  
2727 `_POSIX_C_SOURCE` is defined, the state of `_POSIX_SOURCE` is completely irrelevant.

2728 It is expected that C bindings to future POSIX standards will define new values for  
2729 `_POSIX_C_SOURCE`, with each new value reserving the name space for that new standard, plus  
2730 all earlier POSIX standards. Using a single feature test macro for all standards rather than a  
2731 separate macro for each standard furthers the goal of eventually combining all of the C bindings  
2732 into one standard.

2733 It is further intended that these feature test macros apply only to the headers specified by  
2734 IEEE Std. 1003.1-200x. Implementations are expressly permitted to make visible symbols not  
2735 specified by IEEE Std. 1003.1-200x, within both IEEE Std. 1003.1-200x and other headers, under  
2736 the control of feature test macros that are not defined by IEEE Std. 1003.1-200x.

### 2737 *B.2.2.2 The Name Space*

2738 The reservation of identifiers is paraphrased from the ISO C standard. The text is included  
2739 because it needs to be part of IEEE Std. 1003.1-200x, regardless of possible changes in future  
2740 versions of the ISO C standard.

2741 These identifiers may be used by implementations, particularly for feature test macros.  
2742 Implementations should not use feature test macro names that might be reasonably used by a  
2743 standard.

2744 Including headers more than once is a reasonably common practice, and it should be carried  
2745 forward from the ISO C standard. More significantly, having definitions in more than one  
2746 header is explicitly permitted. Where the potential declaration is “benign” (the same definition  
2747 twice) the declaration can be repeated, if that is permitted by the compiler. (This is usually true  
2748 of macros, for example.) In those situations where a repetition is not benign (for example,

2749 **typedefs**), conditional compilation must be used. The situation actually occurs both within the  
 2750 ISO C standard and within POSIX.1: **time\_t** should be in `<sys/types.h>`, and the ISO C standard  
 2751 mandates that it be in `<time.h>`.

2752 The area of name space pollution *versus* additions to structures is difficult because of the macro  
 2753 structure of C. The following discussion summarizes all the various problems with and  
 2754 objections to the issue.

2755 Note the phrase “user-defined macro”. Users are not permitted to define macro names (or any  
 2756 other name) beginning with “\_`[A-Z_]`”. Thus, the conflict cannot occur for symbols reserved  
 2757 to the vendor’s name space, and the permission to add fields automatically applies, without  
 2758 qualification, to those symbols.

2759 1. Data structures (and unions) need to be defined in headers by implementations to meet  
 2760 certain requirements of POSIX.1 and the ISO C standard.

2761 2. The structures defined by POSIX.1 are typically minimal, and any practical  
 2762 implementation would wish to add fields to these structures either to hold additional  
 2763 related information or for backwards-compatibility (or both). Future standards (and *de*  
 2764 *facto* standards) would also wish to add to these structures. Issues of field alignment make  
 2765 it impractical (at least in the general case) to simply omit fields when they are not defined  
 2766 by the particular standard involved.

2767 Struct **dirent** is an example of such a minimal structure (although one could argue about  
 2768 whether the other fields need visible names). The `st_rdev` field of most implementations’  
 2769 **stat** structure is a common example where extension is needed and where a conflict could  
 2770 occur.

2771 3. Fields in structures are in an independent name space, so the addition of such fields  
 2772 presents no problem to the C language itself in that such names cannot interact with  
 2773 identically named user symbols because access is qualified by the specific structure name.

2774 4. There is an exception to this: macro processing is done at a lexical level. Thus, symbols  
 2775 added to a structure might be recognized as user-provided macro names at the location  
 2776 where the structure is declared. This only can occur if the user-provided name is declared  
 2777 as a macro before the header declaring the structure is included. The user’s use of the name  
 2778 after the declaration cannot interfere with the structure because the symbol is hidden and  
 2779 only accessible through access to the structure. Presumably, the user would not declare  
 2780 such a macro if there was an intention to use that field name.

2781 5. Macros from the same or a related header might use the additional fields in the structure,  
 2782 and those field names might also collide with user macros. Although this is a less frequent  
 2783 occurrence, since macros are expanded at the point of use, no constraint on the order of use  
 2784 of names can apply.

2785 6. An “obvious” solution of using names in the reserved name space and then redefining  
 2786 them as macros when they should be visible does not work because this has the effect of  
 2787 exporting the symbol into the general name space. For example, given a (hypothetical)  
 2788 system-provided header `<h.h>`, and two parts of a C program in `a.c` and `b.c`, in header  
 2789 `<h.h>`:

```

2790     struct foo {
2791         int __i;
2792     }
2793
2794     #ifdef _FEATURE_TEST
2795     #define i __i;
2796 #endif

```

2796 **In file a.c:**

```

2797     #include h.h
2798     extern int i;
2799     ...

```

2800 **In file b.c:**

```

2801     extern int i;
2802     ...

```

2803 The symbol that the user thinks of as *i* in both files has an external name of `__i` in **a.c**; the  
 2804 same symbol *i* in **b.c** has an external name *i* (ignoring any hidden manipulations the  
 2805 compiler might perform on the names). This would cause a mysterious name resolution  
 2806 problem when **a.o** and **b.o** are linked.

2807 Simply avoiding definition then causes alignment problems in the structure.

2808 A structure of the form:

```

2809     struct foo {
2810         union {
2811             int __i;
2812             #ifdef _FEATURE_TEST
2813             int i;
2814             #endif
2815         } __ii;
2816     }

```

2817 does not work because the name of the logical field *i* is `__ii.i`, and introduction of a macro  
 2818 to restore the logical name immediately reintroduces the problem discussed previously  
 2819 (although its manifestation might be more immediate because a syntax error would result  
 2820 if a recursive macro did not cause it to fail first).

2821 7. A more workable solution would be to declare the structure:

```

2822     struct foo {
2823         #ifdef _FEATURE_TEST
2824             int i;
2825         #else
2826             int __i;
2827         #endif
2828     }

```

2829 However, if a macro (particularly one required by a standard) is to be defined that uses  
 2830 this field, two must be defined: one that uses *i*, the other that uses `__i`. If more than one  
 2831 additional field is used in a macro and they are conditional on distinct combinations of  
 2832 features, the complexity goes up as  $2^n$ .

2833 All this leaves a difficult situation: vendors must provide very complex headers to deal with  
 2834 what is conceptually simple and safe—adding a field to a structure. It is the possibility of user-

2835 provided macros with the same name that makes this difficult.

2836 Several alternatives were proposed that involved constraining the user's access to part of the  
2837 name space available to the user (as specified by the ISO C standard). In some cases, this was  
2838 only until all the headers had been included. There were two proposals discussed that failed to  
2839 achieve consensus:

2840 1. Limiting it for the whole program.

2841 2. Restricting the use of identifiers containing only uppercase letters until after all system  
2842 headers had been included. It was also pointed out that because macros might wish to  
2843 access fields of a structure (and macro expansion occurs totally at point of use) restricting  
2844 names in this way would not protect the macro expansion, and thus the solution was  
2845 inadequate.

2846 It was finally decided that reservation of symbols would occur, but as constrained.

2847 The current wording also allows the addition of fields to a structure, but requires that user  
2848 macros of the same name not interfere. This allows vendors to do one of the following:

2849 • Not create the situation (do not extend the structures with user-accessible names or use the  
2850 solution in (7) above)

2851 • Extend their compilers to allow some way of adding names to structures and macros safely

2852 There are at least two ways that the compiler might be extended: add new preprocessor  
2853 directives that turn off and on macro expansion for certain symbols (without changing the value  
2854 of the macro) and a function or lexical operation that suppresses expansion of a word. The latter  
2855 seems more flexible, particularly because it addresses the problem in macros as well as in  
2856 declarations.

2857 The following seems to be a possible implementation extension to the C language that will do  
2858 this: any token that during macro expansion is found to be preceded by three ' #' symbols shall  
2859 not be further expanded in exactly the same way as described for macros that expand to their  
2860 own name as in Section 3.8.3.4 of the ISO C standard. A vendor may also wish to implement this  
2861 as an operation that is lexically a function, which might be implemented as:

```
2862 #define __safe_name(x) ###x
```

2863 Using a function notation would insulate vendors from changes in standards until such a  
2864 functionality is standardized (if ever). Standardization of such a function would be valuable  
2865 because it would then permit third parties to take advantage of it portably in software they may  
2866 supply.

2867 The symbols that are "explicitly permitted, but not required by IEEE Std. 1003.1-200x" include  
2868 those classified below. (That is, the symbols classified below might, but are not required to, be  
2869 present when `_POSIX_C_SOURCE` is defined to have the value 20010xL.)

2870 • Symbols in `<limits.h>` and `<unistd.h>` that are defined to indicate support for options or  
2871 limits that are constant at compile-time.

2872 • Symbols in the name space reserved for the implementation by the ISO C standard.

2873 • Symbols in a name space reserved for a particular type of extension (for example, type names  
2874 ending with `_t` in `<sys/types.h>`).

2875 • Additional members of structures or unions whose names do not reduce the name space  
2876 reserved for applications.

2877 Since both implementations and future revisions of IEEE Std. 1003.1-200x and other POSIX  
2878 standards may use symbols in the reserved spaces described in these tables, there is a potential

2879 for name space clashes. To avoid future name space clashes when adding symbols,  
2880 implementations should not use the `posix_`, `POSIX_`, or `_POSIX_` prefixes.

### 2881 **B.2.3 Error Numbers**

2882 It was the consensus of the standard developers that to allow the conformance document to  
2883 state that an error occurs and under what conditions, but to disallow a statement that it never  
2884 occurs, does not make sense. It could be implied by the current wording that this is allowed, but  
2885 to reduce the possibility of future interpretation requests, it is better to make an explicit  
2886 statement.

2887 The ISO C standard requires that `errno` be an assignable *lvalue*. Originally, the definition in  
2888 POSIX.1 was stricter than that in the ISO C standard, `extern int errno`, in order to support  
2889 historical usage. In a multi-threaded environment, implementing `errno` as a global variable  
2890 results in non-deterministic results when accessed. It is required, however, that `errno` work as a  
2891 per-thread error reporting mechanism. In order to do this, a separate `errno` value has to be  
2892 maintained for each thread. The following section discusses the various alternative solutions  
2893 that were considered.

2894 In order to avoid this problem altogether for new functions, these functions avoid using `errno`  
2895 and, instead, return the error number directly as the function return value; a return value of zero  
2896 indicates that no error was detected.

2897 For any function that can return errors, the function return value is not used for any purpose  
2898 other than for reporting errors. Even when the output of the function is scalar, it is passed  
2899 through a function argument. While it might have been possible to allow some scalar outputs to  
2900 be coded as negative function return values and mixed in with positive error status returns, this  
2901 was rejected—using the return value for a mixed purpose was judged to be of limited use and  
2902 error prone.

2903 Checking the value of `errno` alone is not sufficient to determine the existence or type of an error,  
2904 since it is not required that a successful function call clear `errno`. The variable `errno` should only  
2905 be examined when the return value of a function indicates that the value of `errno` is meaningful.  
2906 In that case, the function is required to set the variable to something other than zero.

2907 The variable `errno` shall never be set to zero by any function call; to do so would contradict the  
2908 ISO C standard.

2909 POSIX.1 requires (in the ERRORS sections of function descriptions) certain error values to be set  
2910 in certain conditions because many existing applications depend on them. Some error numbers,  
2911 such as [EFAULT], are entirely implementation-defined and are noted as such in their  
2912 description in the ERRORS section. This section otherwise allows wide latitude to the  
2913 implementation in handling error reporting.

2914 Some of the ERRORS sections in IEEE Std. 1003.1-200x have two subsections. The first:

2915        “The function shall fail if:”

2916 could be called the “mandatory” section.

2917 The second:

2918        “The function may fail if:”

2919 could be informally known as the “optional” section.

2920 Attempting to infer the quality of an implementation based on whether it detects optional error  
2921 conditions is not useful.



2922		Following each one-word symbolic name for an error, there is a description of the error. The
2923		rationale for some of the symbolic names follows:
2924	[ECANCELED]	This spelling was chosen as being more common.
2925	[EFAULT]	Most historical implementations do not catch an error and set <i>errno</i> when an
2926		invalid address is given to the functions <i>wait()</i> , <i>time()</i> , or <i>times()</i> . Some
2927		implementations cannot reliably detect an invalid address. And most systems
2928		that detect invalid addresses will do so only for a system call, not for a library
2929		routine.
2930	[EFTYPE]	This error code was proposed in earlier proposals as “Inappropriate operation
2931		for file type”, meaning that the operation requested is not appropriate for the
2932		file specified in the function call. This code was proposed, although the same
2933		idea was covered by [ENOTTY], because the connotations of the name would
2934		be misleading. It was pointed out that the <i>fcntl()</i> function uses the error code
2935		[EINVAL] for this notion, and hence all instances of [EFTYPE] were changed
2936		to this code.
2937	[EINTR]	POSIX.1 prohibits conforming implementations from restarting interrupted
2938		system calls. However, it does not require that [EINTR] be returned when
2939		another legitimate value may be substituted; for example, a partial transfer
2940		count when <i>read()</i> or <i>write()</i> are interrupted. This is only given when the
2941		signal catching function returns normally as opposed to returns by
2942		mechanisms like <i>longjmp()</i> or <i>siglongjmp()</i> .
2943	[ELOOP]	In specifying conditions under which implementations would generate this
2944		error, the following goals were considered:
2945		• To ensure that actual loops are detected, including loops that result from
2946		symbolic links across distributed file systems.
2947		• To ensure that during path name resolution an application can rely on the
2948		ability to follow at least {SYMLOOP_MAX} symbolic links in the absence
2949		of a loop.
2950		• To allow implementations to provide the capability of traversing more
2951		than {SYMLOOP_MAX} symbolic links in the absence of a loop.
2952		• To allow implementations to detect loops and generate the error prior to
2953		encountering {SYMLOOP_MAX} symbolic links.
2954	[ENAMETOOLONG]	
2955		When a symbolic link is encountered during path name resolution, the
2956		contents of that symbolic link are used to create a new path name. The
2957		standard developers intended to allow, but not require, that implementations
2958		enforce the restriction of {PATH_MAX} on the result of this path name
2959		substitution.
2960	[ENOMEM]	The term <i>main memory</i> is not used in POSIX.1 because it is implementation-
2961		defined.
2962	[ENOTSUP]	This error code is to be used when an implementation chooses to implement
2963		the required functionality of IEEE Std. 1003.1-200x but does not support
2964		optional facilities defined by IEEE Std. 1003.1-200x. The return of [ENOSYS] is
2965		to be taken to indicate that the function of the interface is not supported at all;
2966		the function will always fail with this error code.

2967 [ENOTTY] The symbolic name for this error is derived from a time when device control  
 2968 was done by *ioctl()* and that operation was only permitted on a terminal  
 2969 interface. The term *TTY* is derived from *teletypewriter*, the devices to which  
 2970 this error originally applied.

2971 [EPIPE] This condition normally generates the signal SIGPIPE; the error is returned if  
 2972 the signal does not terminate the process.

2973 [EROFS] In historical implementations, attempting to *unlink()* or *rmdir()* a mount point  
 2974 would generate an [EBUSY] error. An implementation could be envisioned  
 2975 where such an operation could be performed without error. In this case, if  
 2976 *either* the directory entry or the actual data structures reside on a read-only file  
 2977 system, [EROFS] is the appropriate error to generate. (For example, changing  
 2978 the link count of a file on a read-only file system could not be done, as is  
 2979 required by *unlink()*, and thus an error should be reported.)

2980 Three error numbers, [EDOM], [EILSEQ], and [ERANGE], were added to this section primarily  
 2981 for consistency with the ISO C standard.

### 2982 **Alternative Solutions for Per-Thread *errno***

2983 The usual implementation of *errno* as a single global variable does not work in a multi-threaded  
 2984 environment. In such an environment, a thread may make a POSIX.1 call and get a -1 error  
 2985 return, but before that thread can check the value of *errno*, another thread might have made a  
 2986 second POSIX.1 call that also set *errno*. This behavior is unacceptable in robust programs. There  
 2987 were a number of alternatives that were considered for handling the *errno* problem:

- 2988 • Implement *errno* as a per-thread integer variable.
- 2989 • Implement *errno* as a service that can access the per-thread error number.
- 2990 • Change all POSIX.1 calls to accept an extra status argument and avoid setting *errno*.
- 2991 • Change all POSIX.1 calls to raise a language exception.

2992 The first option offers the highest level of compatibility with existing practice but requires  
 2993 special support in the linker, compiler, and/or virtual memory system to support the new  
 2994 concept of thread private variables. When compared with current practice, the third and fourth  
 2995 options are much cleaner, more efficient, and encourage a more robust programming style, but  
 2996 they require new versions of all of the POSIX.1 functions that might detect an error. The second  
 2997 option offers compatibility with existing code that uses the `<errno.h>` header to define the  
 2998 symbol *errno*. In this option, *errno* may be a macro defined:

```
2999     #define errno  (*__errno())
3000     extern int    *__errno();
```

3001 This option may be implemented as a per-thread variable whereby an *errno* field is allocated in  
 3002 the user space object representing a thread, and whereby the function *\_\_errno()* makes a system  
 3003 call to determine the location of its user space object and returns the address of the *errno* field of  
 3004 that object. Another implementation, one that avoids calling the kernel, involves allocating  
 3005 stacks in chunks. The stack allocator keeps a side table indexed by chunk number containing a  
 3006 pointer to the thread object that uses that chunk. The *\_\_errno()* function then looks at the stack  
 3007 pointer, determines the chunk number, and uses that as an index into the chunk table to find its  
 3008 thread object and thus its private value of *errno*. On most architectures, this can be done in four  
 3009 to five instructions. Some compilers may wish to implement *\_\_errno()* inline to improve  
 3010 performance.

**3011 Disallowing Return of the [EINTR] Error Code**

3012 Many blocking interfaces defined by IEEE Std. 1003.1-200x may return [EINTR] if interrupted  
3013 during their execution by a signal handler. Blocking interfaces introduced under the Threads  
3014 option do not have this property. Instead, they require that the interface appear to be atomic  
3015 with respect to interruption. In particular, clients of block interfaces need not handle any  
3016 possible [EINTR] return as a special case since it will never occur. If it is necessary to restart  
3017 operations or complete incomplete operations following the execution of a signal handler, this is  
3018 handled by the implementation, rather than by the application.

3019 Requiring applications to handle [EINTR] errors on blocking interfaces has been shown to be a  
3020 frequent source of often unreproducible bugs, and it adds no compelling value to the available  
3021 functionality. Thus, blocking interfaces introduced for use by multi-threaded programs do not  
3022 use this paradigm. In particular, in none of the functions *flockfile()*, *pthread\_cond\_timedwait()*,  
3023 *pthread\_cond\_wait()*, *pthread\_join()*, *pthread\_mutex\_lock()*, and *sigwait()* did providing [EINTR]  
3024 returns add value, or even particularly make sense. Thus, these functions do not provide for an  
3025 [EINTR] return, even when interrupted by a signal handler. The same arguments can be applied  
3026 to *sem\_wait()*, *sem\_trywait()*, *sigwaitinfo()*, and *sigtimedwait()*, but implementations are  
3027 permitted to return [EINTR] error codes for these functions for compatibility with earlier  
3028 versions of IEEE Std. 1003.1-200x. Applications cannot rely on calls to these functions returning  
3029 [EINTR] error codes when signals are delivered to the calling thread, but they should allow for  
3030 the possibility.

**3031 B.2.3.1 Additional Error Numbers**

3032 The ISO C standard defines the name space for implementations to add additional error  
3033 numbers.

**3034 B.2.4 Signal Concepts**

3035 Historical implementations of signals, using the *signal()* function, have shortcomings that make  
3036 them unreliable for many application uses. Because of this, a new signal mechanism, based very  
3037 closely on the one of 4.2 BSD and 4.3 BSD, was added to POSIX.1.

**3038 Signal Names**

3039 The restriction on the actual type used for **sigset\_t** is intended to guarantee that these objects can  
3040 always be assigned, have their address taken, and be passed as parameters by value. It is not  
3041 intended that this type be a structure including pointers to other data structures, as that could  
3042 impact the portability of applications performing such operations. A reasonable implementation  
3043 could be a structure containing an array of some integer type.

3044 The signals described in IEEE Std. 1003.1-200x must have unique values so that they may be  
3045 named as parameters of **case** statements in the body of a C language **switch** clause. However,  
3046 implementation-defined signals may have values that overlap with each other or with signals  
3047 specified in IEEE Std. 1003.1-200x. An example of this is SIGABRT, which traditionally overlaps  
3048 some other signal, such as SIGIOT.

3049 SIGKILL, SIGTERM, SIGUSR1, and SIGUSR2 are ordinarily generated only through the explicit  
3050 use of the *kill()* function, although some implementations generate SIGKILL under  
3051 extraordinary circumstances. SIGTERM is traditionally the default signal sent by the *kill*  
3052 command.

3053 The signals SIGBUS, SIGEMT, SIGIOT, SIGTRAP, and SIGSYS were omitted from POSIX.1  
3054 because their behavior is implementation-defined and could not be adequately categorized.  
3055 Conforming implementations may deliver these signals, but must document the circumstances

3056 under which they are delivered and note any restrictions concerning their delivery. The signals  
3057 SIGFPE, SIGILL, and SIGSEGV are similar in that they also generally result only from  
3058 programming errors. They were included in POSIX.1 because they do indicate three relatively  
3059 well-categorized conditions. They are all defined by the ISO C standard and thus would have to  
3060 be defined by any system with a ISO C standard binding, even if not explicitly included in  
3061 POSIX.1.

3062 There is very little that a Conforming POSIX.1 Application can do by catching, ignoring, or  
3063 masking any of the signals SIGILL, SIGTRAP, SIGIOT, SIGEMT, SIGBUS, SIGSEGV, SIGSYS, or  
3064 SIGFPE. They will generally be generated by the system only in cases of programming errors.  
3065 While it may be desirable for some robust code (for example, a library routine) to be able to  
3066 detect and recover from programming errors in other code, these signals are not nearly sufficient  
3067 for that purpose. One portable use that does exist for these signals is that a command interpreter  
3068 can recognize them as the cause of a process' termination (with *wait()*) and print an appropriate  
3069 message. The mnemonic tags for these signals are derived from their PDP-11 origin.

3070 The signals SIGSTOP, SIGTSTP, SIGTTIN, SIGTTOU, and SIGCONT are provided for job control  
3071 and are unchanged from 4.2 BSD. The signal SIGCHLD is also typically used by job control  
3072 shells to detect children that have terminated or, as in 4.2 BSD, stopped.

3073 Some implementations, including System V, have a signal named SIGCLD, which is similar to  
3074 SIGCHLD in 4.2 BSD. POSIX.1 permits implementations to have a single signal with both  
3075 names. POSIX.1 carefully specifies ways in which portable applications can avoid the semantic  
3076 differences between the two different implementations. The name SIGCHLD was chosen for  
3077 POSIX.1 because most current application usages of it can remain unchanged in conforming  
3078 applications. SIGCLD in System V has more cases of semantics that POSIX.1 does not specify,  
3079 and thus applications using it are more likely to require changes in addition to the name change.

3080 The signals SIGUSR1 and SIGUSR2 are commonly used by applications for notification of  
3081 exceptional behavior and are described as "reserved as application-defined" so that such use is  
3082 not prohibited. Implementations should not generate SIGUSR1 or SIGUSR2, except when  
3083 explicitly requested by *kill()*. It is recommended that libraries not use these two signals, as such  
3084 use in libraries could interfere with their use by applications calling the libraries. If such use is  
3085 unavoidable, it should be documented. It is prudent for non-portable libraries to use non-  
3086 standard signals to avoid conflicts with use of standard signals by portable libraries.

3087 There is no portable way for an application to catch or ignore non-standard signals. Some  
3088 implementations define the range of signal numbers, so applications can install signal-catching  
3089 functions for all of them. Unfortunately, implementation-defined signals often cause problems  
3090 when caught or ignored by applications that do not understand the reason for the signal. While  
3091 the desire exists for an application to be more robust by handling all possible signals (even those  
3092 only generated by *kill()*), no existing mechanism was found to be sufficiently portable to include  
3093 in POSIX.1. The value of such a mechanism, if included, would be diminished given that  
3094 SIGKILL would still not be catchable.

3095 A number of new signal numbers are reserved for applications because the two user signals  
3096 defined by POSIX.1 are insufficient for many realtime applications. A range of signal numbers is  
3097 specified, rather than an enumeration of additional reserved signal names, because different  
3098 applications and application profiles will require a different number of application signals. It is  
3099 not desirable to burden all application domains and therefore all implementations with the  
3100 maximum number of signals required by all possible applications. Note that in this context,  
3101 signal numbers are essentially different signal priorities.

3102 The relatively small number of required additional signals, `{_POSIX_RTSIG_MAX}`, was chosen  
3103 so as not to require an unreasonably large signal mask/set. While this number of signals defined  
3104 in POSIX.1 will fit in a single 32-bit word signal mask, it is recognized that most existing

3105 implementations define many more signals than are specified in POSIX.1 and, in fact, many  
3106 implementations have already exceeded 32 signals (including the “null signal”). Support of  
3107 `{_POSIX_RTSIG_MAX}` additional signals may push some implementation over the single 32-bit  
3108 word line, but is unlikely to push any implementations that are already over that line beyond the  
3109 64-signal line.

#### 3110 *B.2.4.1 Signal Generation and Delivery*

3111 The terms defined in this section are not used consistently in documentation of historical  
3112 systems. Each signal can be considered to have a lifetime beginning with *generation* and ending  
3113 with *delivery* or *acceptance*. The POSIX.1 definition of *delivery* does not exclude ignored signals;  
3114 this is considered a more consistent definition. This revised text in several parts of  
3115 IEEE Std. 1003.1-200x clarifies the distinct semantics of asynchronous signal *delivery* and  
3116 synchronous signal *acceptance*. The previous wording attempted to categorize both under the  
3117 term *delivery*, which led to conflicts over whether the effects of asynchronous signal delivery  
3118 applied to synchronous signal acceptance.

3119 Signals generated for a process are delivered to only one thread. Thus, if more than one thread is  
3120 eligible to receive a signal, one has to be chosen. The choice of threads is left entirely up to the  
3121 implementation both to allow the widest possible range of conforming implementations and to  
3122 give implementations the freedom to deliver the signal to the “easiest possible” thread should  
3123 there be differences in ease of delivery between different threads.

3124 Note that should multiple delivery among cooperating threads be required by an application,  
3125 this can be trivially constructed out of the provided single-delivery semantics. The construction  
3126 of a *sigwait\_multiple()* function that accomplishes this goal is presented with the rationale for  
3127 *sigwaitinfo()*.

3128 Implementations should deliver unblocked signals as soon after they are generated as possible.  
3129 However, it is difficult for POSIX.1 to make specific requirements about this, beyond those in  
3130 *kill()* and *sigprocmask()*. Even on systems with prompt delivery, scheduling of higher priority  
3131 processes is always likely to cause delays.

3132 In general, the interval between the generation and delivery of unblocked signals cannot be  
3133 detected by an application. Thus, references to pending signals generally apply to blocked,  
3134 pending signals. An implementation registers a signal as pending on the process when no thread  
3135 has the signal unblocked and there are no threads blocked in a *sigwait()* function for that signal.  
3136 Thereafter, the implementation delivers the signal to the first thread that unblocks the signal or  
3137 calls a *sigwait()* function on a signal set containing this signal rather than choosing the recipient  
3138 thread at the time the signal is sent.

3139 In the 4.3 BSD system, signals that are blocked and set to `SIG_IGN` are discarded immediately  
3140 upon generation. For a signal that is ignored as its default action, if the action is `SIG_DFL` and  
3141 the signal is blocked, a generated signal remains pending. In the 4.1 BSD system and in  
3142 System V, Release 3, two other implementations that support a somewhat similar signal  
3143 mechanism, all ignored, blocked signals remain pending if generated. Because it is not normally  
3144 useful for an application to simultaneously ignore and block the same signal, it was unnecessary  
3145 for POSIX.1 to specify behavior that would invalidate any of the historical implementations.

3146 There is one case in some historical implementations where an unblocked, pending signal does  
3147 not remain pending until it is delivered. In the System V implementation of *signal()*, pending  
3148 signals are discarded when the action is set to `SIG_DFL` or a signal-catching routine (as well as to  
3149 `SIG_IGN`). Except in the case of setting `SIGCHLD` to `SIG_DFL`, implementations that do this do  
3150 not conform completely to POSIX.1. Some earlier proposals for POSIX.1 explicitly stated this,  
3151 but these statements were redundant due to the requirement that functions defined by POSIX.1  
3152 not change attributes of processes defined by POSIX.1 except as explicitly stated.

3153 POSIX.1 specifically states that the order in which multiple, simultaneously pending signals are  
 3154 delivered is unspecified. This order has not been explicitly specified in historical  
 3155 implementations, but has remained quite consistent and been known to those familiar with the  
 3156 implementations. Thus, there have been cases where applications (usually system utilities) have  
 3157 been written with explicit or implicit dependencies on this order. Implementors and others  
 3158 porting existing applications may need to be aware of such dependencies.

3159 When there are multiple pending signals that are not blocked, implementations should arrange  
 3160 for the delivery of all signals at once, if possible. Some implementations stack calls to all pending  
 3161 signal-catching routines, making it appear that each signal-catcher was interrupted by the next  
 3162 signal. In this case, the implementation should ensure that this stacking of signals does not  
 3163 violate the semantics of the signal masks established by *sigaction()*. Other implementations  
 3164 process at most one signal when the operating system is entered, with remaining signals saved  
 3165 for later delivery. Although this practice is widespread, this behavior is neither standardized  
 3166 nor endorsed. In either case, implementations should attempt to deliver signals associated with  
 3167 the current state of the process (for example, SIGFPE) before other signals, if possible.

3168 In 4.2 BSD and 4.3 BSD, it is not permissible to ignore or explicitly block SIGCONT, because if  
 3169 blocking or ignoring this signal prevented it from continuing a stopped process, such a process  
 3170 could never be continued (only killed by SIGKILL). However, 4.2 BSD and 4.3 BSD do block  
 3171 SIGCONT during execution of its signal-catching function when it is caught, creating exactly  
 3172 this problem. A proposal was considered to disallow catching SIGCONT in addition to ignoring  
 3173 and blocking it, but this limitation led to objections. The consensus was to require that  
 3174 SIGCONT always continue a stopped process when generated. This removed the need to  
 3175 disallow ignoring or explicit blocking of the signal; note that SIG\_IGN and SIG\_DFL are  
 3176 equivalent for SIGCONT .

#### 3177 B.2.4.2 Realtime Signal Generation and Delivery

3178 The Realtime Signals Extension option to POSIX.1 signal generation and delivery behavior is  
 3179 required for the following reasons:

- 3180 • The **sigevent** structure is used by other POSIX.1 functions that result in asynchronous event  
 3181 notifications to specify the notification mechanism to use and other information needed by  
 3182 the notification mechanism. IEEE Std. 1003.1-200x defines only three symbolic values for the  
 3183 notification mechanism. SIGEV\_NONE is used to indicate that no notification is required  
 3184 when the event occurs. This is useful for applications that use asynchronous I/O with polling  
 3185 for completion. SIGEV\_SIGNAL indicates that a signal shall be generated when the event  
 3186 occurs. SIGEV\_NOTIFY provides for “callback functions” for asynchronous notifications  
 3187 done by a function call within the context of a new thread. This provides a multi-threaded  
 3188 process a more natural means of notification than signals. The primary difficulty with  
 3189 previous notification approaches has been to specify the environment of the notification  
 3190 routine.

- 3191 — One approach is to limit the notification routine to call only functions permitted in a  
 3192 signal handler. While the list of permissible functions is clearly stated, this is overly  
 3193 restrictive.

- 3194 — A second approach is to define a new list of functions or classes of functions that are  
 3195 explicitly permitted or not permitted. This would give a programmer more lists to deal  
 3196 with, which would be awkward.

- 3197 — The third approach is to define completely the environment for execution of the  
 3198 notification function. A clear definition of an execution environment for notification is  
 3199 provided by executing the notification function in the environment of a newly created  
 3200 thread.

- 3201 Implementations may support additional notification mechanisms by defining new values  
3202 for *sigev\_notify*.
- 3203 For a notification type of SIGEV\_SIGNAL, the other members of the **sigevent** structure  
3204 defined by IEEE Std. 1003.1-200x specify the realtime signal—that is, the signal number and  
3205 application-defined value that differentiates between occurrences of signals with the same  
3206 number—that will be generated when the event occurs. The structure is defined in  
3207 <**signal.h**>, even though the structure is not directly used by any of the signal functions,  
3208 because it is part of the signals interface used by the POSIX.1b “client functions”. When the  
3209 client functions include <**signal.h**> to define the signal names, the **sigevent** structure will  
3210 also be defined.
- 3211 An application-defined value passed to the signal handler is used to differentiate between  
3212 different “events” instead of requiring that the application use different signal numbers for  
3213 several reasons:
- 3214 — Realtime applications potentially handle a very large number of different events.  
3215 Requiring that implementations support a correspondingly large number of distinct  
3216 signal numbers will adversely impact the performance of signal delivery because the  
3217 signal masks to be manipulated on entry and exit to the handlers will become large.
  - 3218 — Event notifications are prioritized by signal number (the rationale for this is explained in  
3219 the following paragraphs) and the use of different signal numbers to differentiate  
3220 between the different event notifications overloads the signal number more than has  
3221 already been done. It also requires that the application writer make arbitrary assignments  
3222 of priority to events that are logically of equal priority.
- 3223 A union is defined for the application-defined value so that either an integer constant or a  
3224 pointer can be portably passed to the signal-catching function. On some architectures a  
3225 pointer cannot be cast to an **int** and *vice versa*.
- 3226 Use of a structure here with an explicit notification type discriminant rather than explicit  
3227 parameters to realtime functions, or embedded in other realtime structures, provides for  
3228 future extensions to IEEE Std. 1003.1-200x. Additional, perhaps more efficient, notification  
3229 mechanisms can be supported for existing realtime function interfaces, such as timers and  
3230 asynchronous I/O, by extending the **sigevent** structure appropriately. The existing realtime  
3231 function interfaces will not have to be modified to use any such new notification mechanism.  
3232 The revised text concerning the SIGEV\_SIGNAL value makes consistent the semantics of the  
3233 members of the **sigevent** structure, particularly in the definitions of *lio\_listio()* and  
3234 *aio\_fsync()*. For uniformity, other revisions cause this specification to be referred to rather  
3235 than inaccurately duplicated in the descriptions of functions and structures using the  
3236 **sigevent** structure. The revised wording does not relax the requirement that the signal  
3237 number be in the range SIGRTMIN to SIGRTMAX to guarantee queuing and passing of the  
3238 application value, since that requirement is still implied by the signal names.
- 3239 • IEEE Std. 1003.1-200x is intentionally vague on whether “non-realtime” signal-generating  
3240 mechanisms can result in a **siginfo\_t** being supplied to the handler on delivery. In one  
3241 existing implementation, a **siginfo\_t** is posted on signal generation, even though the  
3242 implementation does not support queuing of multiple occurrences of a signal. It is not the  
3243 intent of IEEE Std. 1003.1-200x to preclude this, independent of the mandate to define signals  
3244 that do support queuing. Any interpretation that appears to preclude this is a mistake in the  
3245 reading or writing of the standard.
  - 3246 • Signals handled by realtime signal handlers might be generated by functions or conditions  
3247 that do not allow the specification of an application-defined value and do not queue.  
3248 IEEE Std. 1003.1-200x specifies the *si\_code* member of the **siginfo\_t** structure used in existing

3249 practice and defines additional codes so that applications can detect whether an application-  
 3250 defined value is present or not. The code SI\_USER for *kill()*-generated signals is adopted  
 3251 from existing practice.

3252 • The *sigaction()* *sa\_flags* value SA\_SIGINFO tells the implementation that the signal-catching  
 3253 function expects two additional arguments. When the flag is not set, a single argument, the  
 3254 signal number, is passed as specified by IEEE Std. 1003.1-200x. Although  
 3255 IEEE Std. 1003.1-200x does not explicitly allow the *info* argument to the handler function to  
 3256 be NULL, this is existing practice. This provides for compatibility with programs whose  
 3257 signal-catching functions are not prepared to accept the additional arguments.  
 3258 IEEE Std. 1003.1-200x is explicitly unspecified as to whether signals actually queue when  
 3259 SA\_SIGINFO is not set for a signal, as there appear to be no benefits to applications in  
 3260 specifying one behavior or another. One existing implementation queues a **siginfo\_t** on each  
 3261 signal generation, unless the signal is already pending, in which case the implementation  
 3262 discards the new **siginfo\_t**; that is, the queue length is never greater than one. This  
 3263 implementation only examines SA\_SIGINFO on signal delivery, discarding the queued  
 3264 **siginfo\_t** if its delivery was not requested.

3265 IEEE Std. 1003.1-200x specifies several new values for the *si\_code* member of the **siginfo\_t**  
 3266 structure. In existing practice, a *si\_code* value of less than or equal to zero indicates that  
 3267 the signal was generated by a process via the *kill()* function. In existing practice, values of *si\_code*  
 3268 that provide additional information for implementation-generated signals, such as SIGFPE or  
 3269 SIGSEGV, are all positive. Thus, if implementations define the new constants specified in  
 3270 IEEE Std. 1003.1-200x to be negative numbers, programs written to use existing practice will  
 3271 not break. IEEE Std. 1003.1-200x chose not to attempt to specify existing practice values of  
 3272 *si\_code* other than SI\_USER both because it was deemed beyond the scope of  
 3273 IEEE Std. 1003.1-200x and because many of the values in existing practice appear to be  
 3274 platform and implementation-defined. But, IEEE Std. 1003.1-200x does specify that if an  
 3275 implementation—for example, one that does not have existing practice in this area—chooses  
 3276 to define additional values for *si\_code*, these values have to be different from the values of the  
 3277 symbols specified by IEEE Std. 1003.1-200x. This will allow portable applications to  
 3278 differentiate between signals generated by one of the POSIX.1b asynchronous events and  
 3279 those generated by other implementation events in a manner compatible with existing  
 3280 practice.

3281 The unique values of *si\_code* for the POSIX.1b asynchronous events have implications for  
 3282 implementations of, for example, asynchronous I/O or message passing in user space library  
 3283 code. Such an implementation will be required to provide a hidden interface to the signal  
 3284 generation mechanism that allows the library to specify the standard values of *si\_code*.

3285 Existing practice also defines additional members of **siginfo\_t**, such as the process ID and  
 3286 user ID of the sending process for *kill()*-generated signals. These members were deemed not  
 3287 necessary to meet the requirements of realtime applications and are not specified by  
 3288 IEEE Std. 1003.1-200x. Neither are they precluded.

3289 The third argument to the signal-catching function, *context*, is left undefined by  
 3290 IEEE Std. 1003.1-200x, but is specified in the interface because it matches existing practice for  
 3291 the SA\_SIGINFO flag. It was considered undesirable to require a separate implementation  
 3292 for SA\_SIGINFO for POSIX conformance on implementations that already support the two  
 3293 additional parameters.

3294 • The requirement to deliver lower numbered signals in the range SIGRTMIN to SIGRTMAX  
 3295 first, when multiple unblocked signals are pending, results from several considerations:

3296 — A method is required to prioritize event notifications. The signal number was chosen  
 3297 instead of, for instance, associating a separate priority with each request, because an



3298 implementation has to check pending signals at various points and select one for delivery  
 3299 when more than one is pending. Specifying a selection order is the minimal additional  
 3300 semantic that will achieve prioritized delivery. If a separate priority were to be associated  
 3301 with queued signals, it would be necessary for an implementation to search all non-  
 3302 empty, non-blocked signal queues and select from among them the pending signal with  
 3303 the highest priority. This would significantly increase the cost of and decrease the  
 3304 determinism of signal delivery.

3305 — Given the specified selection of the lowest numeric unblocked pending signal,  
 3306 preemptive priority signal delivery can be achieved using signal numbers and signal  
 3307 masks by ensuring that the *sa\_mask* for each signal number blocks all signals with a  
 3308 higher numeric value.

3309 For realtime applications that want to use only the newly defined realtime signal numbers  
 3310 without interference from the standard signals, this can be achieved by blocking all of the  
 3311 standard signals in the process signal mask and in the *sa\_mask* installed by the signal  
 3312 action for the realtime signal handlers.

3313 IEEE Std. 1003.1-200x explicitly leaves unspecified the ordering of signals outside of the  
 3314 range of realtime signals and the ordering of signals within this range with respect to those  
 3315 outside the range. It was believed that this would unduly constrain implementations or  
 3316 standards in the future definition of new signals.

### 3317 B.2.4.3 Signal Actions

3318 Early proposals mentioned SIGCONT as a second exception to the rule that signals are not  
 3319 delivered to stopped processes until continued. Because IEEE Std. 1003.1-200x now specifies that  
 3320 SIGCONT causes the stopped process to continue when it is generated, delivery of SIGCONT is  
 3321 not prevented because a process is stopped, even without an explicit exception to this rule.

3322 Ignoring a signal by setting the action to SIG\_IGN (or SIG\_DFL for signals whose default action  
 3323 is to ignore) is not the same as installing a signal-catching function that simply returns. Invoking  
 3324 such a function will interrupt certain system functions that block processes (for example, *wait()*,  
 3325 *sigsuspend()*, *pause()*, *read()*, *write()*) while ignoring a signal has no such effect on the process.

3326 Historical implementations discard pending signals when the action is set to SIG\_IGN.  
 3327 However, they do not always do the same when the action is set to SIG\_DFL and the default  
 3328 action is to ignore the signal. IEEE Std. 1003.1-200x requires this for the sake of consistency and  
 3329 also for completeness, since the only signal this applies to is SIGCHLD, and  
 3330 IEEE Std. 1003.1-200x disallows setting its action to SIG\_IGN.

3331 The specification of the effects of SIG\_IGN on SIGCHLD as implementation-defined permits,  
 3332 but does not require, the System V effect of causing terminating children to be ignored by *wait()*.  
 3333 Yet it permits SIGCHLD to be effectively ignored in an implementation-defined manner by use  
 3334 of SIG\_DFL.

3335 Some implementations (System V, for example) assign different semantics for SIGCLD  
 3336 depending on whether the action is set to SIG\_IGN or SIG\_DFL. Since POSIX.1 requires that the  
 3337 default action for SIGCHLD be to ignore the signal, applications should always set the action to  
 3338 SIG\_DFL in order to avoid SIGCHLD.

3339 Some implementations (System V, for example) will deliver a SIGCLD signal immediately when  
 3340 a process establishes a signal-catching function for SIGCLD when that process has a child that  
 3341 has already terminated. Other implementations, such as 4.3 BSD, do not generate a new  
 3342 SIGCHLD signal in this way. In general, a process should not attempt to alter the signal action  
 3343 for the SIGCHLD signal while it has any outstanding children. However, it is not always  
 3344 possible for a process to avoid this; for example, shells sometimes start up processes in pipelines

3345 with other processes from the pipeline as children. Processes that cannot ensure that they have  
3346 no children when altering the signal action for SIGCHLD thus need to be prepared for, but not  
3347 depend on, generation of an immediate SIGCHLD signal.

3348 The default action of the stop signals (SIGSTOP , SIGTSTP, SIGTTIN, SIGTTOU) is to stop a  
3349 process that is executing. If a stop signal is delivered to a process that is already stopped, it has  
3350 no effect. In fact, if a stop signal is generated for a stopped process whose signal mask blocks the  
3351 signal, the signal will never be delivered to the process since the process must receive a  
3352 SIGCONT, which discards all pending stop signals, in order to continue executing.

3353 The SIGCONT signal shall continue a stopped process even if SIGCONT is blocked (or ignored).  
3354 However, if a signal-catching routine has been established for SIGCONT, it will not be entered  
3355 until SIGCONT is unblocked.

3356 If a process in an orphaned process group stops, it is no longer under the control of a job control  
3357 shell and hence would not normally ever be continued. Because of this, orphaned processes that  
3358 receive terminal-related stop signals (SIGTSTP , SIGTTIN, SIGTTOU, but not SIGSTOP ) must  
3359 not be allowed to stop. The goal is to prevent stopped processes from languishing forever. (As  
3360 SIGSTOP is sent only via *kill()*, it is assumed that the process or user sending a SIGSTOP can  
3361 send a SIGCONT when desired.) Instead, the system must discard the stop signal. As an  
3362 extension, it may also deliver another signal in its place. 4.3 BSD sends a SIGKILL, which is  
3363 overly effective because SIGKILL is not catchable. Another possible choice is SIGHUP. 4.3 BSD  
3364 also does this for orphaned processes (processes whose parent has terminated) rather than for  
3365 members of orphaned process groups; this is less desirable because job control shells manage  
3366 process groups. POSIX.1 also prevents SIGTTIN and SIGTTOU signals from being generated for  
3367 processes in orphaned process groups as a direct result of activity on a terminal, preventing  
3368 infinite loops when *read()* and *write()* calls generate signals that are discarded; see Section  
3369 A.11.1.4 (on page 3371). A similar restriction on the generation of SIGTSTP was considered, but  
3370 that would be unnecessary and more difficult to implement due to its asynchronous nature.

3371 Although POSIX.1 requires that signal-catching functions be called with only one argument,  
3372 there is nothing to prevent conforming implementations from extending POSIX.1 to pass  
3373 additional arguments, as long as Strictly Conforming POSIX.1 Applications continue to compile  
3374 and execute correctly. Most historical implementations do, in fact, pass additional, signal-  
3375 specific arguments to certain signal-catching routines.

3376 There was a proposal to change the declared type of the signal handler to:

```
3377     void func (int sig, ...);
```

3378 The usage of ellipses ("...") is ISO C standard syntax to indicate a variable number of  
3379 arguments. Its use was intended to allow the implementation to pass additional information to  
3380 the signal handler in a standard manner.

3381 Unfortunately, this construct would require all signal handlers to be defined with this syntax  
3382 because the ISO C standard allows implementations to use a different parameter passing  
3383 mechanism for variable parameter lists than for non-variable parameter lists. Thus, all existing  
3384 signal handlers in all existing applications would have to be changed to use the variable syntax  
3385 in order to be standard and portable. This is in conflict with the goal of Minimal Changes to  
3386 Existing Application Code.

3387 When terminating a process from a signal-catching function, processes should be aware of any  
3388 interpretation that their parent may make of the status returned by *wait()* or *waitpid()*. In  
3389 particular, a signal-catching function should not call *exit(0)* or *\_exit(0)* unless it wants to indicate  
3390 successful termination. A non-zero argument to *exit()* or *\_exit()* can be used to indicate  
3391 unsuccessful termination. Alternatively, the process can use *kill()* to send itself a fatal signal  
3392 (first ensuring that the signal is set to the default action and not blocked). See also the

3393 RATIONALE section of the `_exit()` function.

3394 The behavior of *unsafe* functions, as defined by this section, is undefined when they are invoked  
3395 from signal-catching functions in certain circumstances. The behavior of reentrant functions, as  
3396 defined by this section, is as specified by POSIX.1, regardless of invocation from a signal-  
3397 catching function. This is the only intended meaning of the statement that reentrant functions  
3398 may be used in signal-catching functions without restriction. Applications must still consider all  
3399 effects of such functions on such things as data structures, files, and process state. In particular,  
3400 application writers need to consider the restrictions on interactions when interrupting `sleep()`  
3401 (see `sleep()`) and interactions among multiple handles for a file description. The fact that any  
3402 specific function is listed as reentrant does not necessarily mean that invocation of that function  
3403 from a signal-catching function is recommended.

3404 In order to prevent errors arising from interrupting non-reentrant function calls, applications  
3405 should protect calls to these functions either by blocking the appropriate signals or through the  
3406 use of some programmatic semaphore. POSIX.1 does not address the more general problem of  
3407 synchronizing access to shared data structures. Note in particular that even the “safe” functions  
3408 may modify the global variable `errno`; the signal-catching function may want to save and restore  
3409 its value. The same principles apply to the reentrancy of application routines and asynchronous  
3410 data access.

3411 Note that `longjmp()` and `siglongjmp()` are not in the list of reentrant functions. This is because the  
3412 code executing after `longjmp()` or `siglongjmp()` can call any unsafe functions with the same  
3413 danger as calling those unsafe functions directly from the signal handler. Applications that use  
3414 `longjmp()` or `siglongjmp()` out of signal handlers require rigorous protection in order to be  
3415 portable. Many of the other functions that are excluded from the list are traditionally  
3416 implemented using either the C language `malloc()` or `free()` functions or the ISO C standard I/O  
3417 library, both of which traditionally use data structures in a non-reentrant manner. Because any  
3418 combination of different functions using a common data structure can cause reentrancy  
3419 problems, POSIX.1 does not define the behavior when any unsafe function is called in a signal  
3420 handler that interrupts any unsafe function.

3421 The only realtime extension to signal actions is the addition of the additional parameters to the  
3422 signal-catching function. This extension has been explained and motivated in the previous  
3423 section. In making this extension, though, developers of POSIX.1b ran into issues relating to  
3424 function prototypes. In response to input from the POSIX.1 standard developers, members were  
3425 added to the **sigaction** structure to specify function prototypes for the newer signal-catching  
3426 function specified by POSIX.1b. These members follow changes that are being made to POSIX.1.  
3427 Note that IEEE Std. 1003.1-200x explicitly states that these fields may overlap so that a union can  
3428 be defined. This will enable existing implementations of POSIX.1 to maintain binary-  
3429 compatibility when these extensions are added.

3430 The **siginfo\_t** structure was adopted for passing the application-defined value to match existing  
3431 practice, but the existing practice has no provision for an application-defined value, so this was  
3432 added. Note that POSIX normally reserves the “\_t” type designation for opaque types. The  
3433 **siginfo\_t** structure breaks with this convention to follow existing practice and thus promote  
3434 portability. Standardization of the existing practice for the other members of this structure may  
3435 be addressed in the future.

3436 Although it is not explicitly visible to applications, there are additional semantics for signal  
3437 actions implied by queued signals and their interaction with other POSIX.1b realtime functions.  
3438 Specifically:

- 3439 • It is not necessary to queue signals whose action is `SIG_IGN`.

- 3440           • For implementations that support POSIX.1b timers, some interaction with the timer functions  
3441           at signal delivery is implied to manage the timer overrun count.

#### 3442 **B.2.4.4** *Signal Effects on Other Functions*

3443           The most common behavior of an interrupted function after a signal-catching function returns is  
3444           for the interrupted function to give an [EINTR] error. However, there are a number of specific  
3445           exceptions, including *sleep()* and certain situations with *read()* and *write()*.

3446           The historical implementations of many functions defined by IEEE Std. 1003.1-200x are not  
3447           interruptible, but delay delivery of signals generated during their execution until after they  
3448           complete. This is never a problem for functions that are guaranteed to complete in a short  
3449           (imperceptible to a human) period of time. It is normally those functions that can suspend a  
3450           process indefinitely or for long periods of time (for example, *wait()*, *pause()*, *sigsuspend()*, *sleep()*,  
3451           or *read()/write()* on a slow device like a terminal] that are interruptible. This permits  
3452           applications to respond to interactive signals or to set timeouts on calls to most such functions  
3453           with *alarm()*. Therefore, implementations should generally make such functions (including ones  
3454           defined as extensions) interruptible.

3455           Functions not mentioned explicitly as interruptible may be so on some implementations,  
3456           possibly as an extension where the function gives an [EINTR] error. There are several functions  
3457           (for example, *getpid()*, *getuid()*) that are specified as never returning an error, which can thus  
3458           never be extended in this way.

### 3459 **B.2.5** **Standard I/O Streams**

#### 3460 **B.2.5.1** *Interaction of File Descriptors and Standard I/O Streams*

3461           There is no additional rationale for this section.

#### 3462 **B.2.5.2** *Stream Orientation and Encoding Rules*

3463           There is no additional rationale for this section.

### 3464 **B.2.6** **STREAMS**

3465           STREAMS are introduced into IEEE Std. 1003.1-200x as part of the alignment with the Single  
3466           UNIX Specification, but marked as an option in recognition that not all systems may wish to  
3467           implement the facility. The option within IEEE Std. 1003.1-200x is denoted by the XSR margin  
3468           marker. The standard developers made this option independent of the XSI option.

3469           STREAMS are a method of implementing network services and other character-based  
3470           input/output mechanisms, with the STREAM being a full-duplex connection between a process  
3471           and a device. STREAMS provides direct access to protocol modules, and optional protocol  
3472           modules can be interposed between the process-end of the STREAM and the device-driver at the  
3473           device-end of the STREAM. Pipes can be implemented using the STREAMS mechanism, so they  
3474           can provide process-to-process as well as process-to-device communications.

3475           This section introduces STREAMS I/O, the message types used to control them, an overview of  
3476           the priority mechanism, and the interfaces used to access them.

3477 **B.2.6.1** *Accessing STREAMS*

3478 There is no additional rationale for this section.

3479 **B.2.7** **XSI Interprocess Communication**

3480 There are two forms of IPC supported as options in IEEE Std. 1003.1-200x. The traditional  
3481 System V IPC routines derived from the SVID—that is, the *msg\**(), *sem\**(), and *shm\**()  
3482 interfaces—are mandatory on XSI-conformant systems. Thus, all XSI-conformant systems  
3483 provide the same mechanisms for manipulating messages, shared memory, and semaphores.

3484 In addition, the POSIX Realtime Extension provides an alternate set of routines for those systems  
3485 supporting the appropriate options.

3486 For maximum portability to UNIX systems, the former are recommended. However, if the  
3487 target for an application is a realtime system, then application developers are advised to write  
3488 their code in such a way that modules using IPC interfaces can be modified easily in the future  
3489 to use either interfaces.

3490 **B.2.7.1** *IPC General Information*

3491 General information that is shared by all three mechanisms is described in this section. The  
3492 common permissions mechanism is briefly introduced, describing the mode bits, and how they  
3493 are used to determine whether or not a process has access to read or write/alter the appropriate  
3494 instance of one of the IPC mechanisms. All other relevant information is contained in the  
3495 reference pages themselves.

3496 The semaphore type of IPC allows processes to communicate through the exchange of  
3497 semaphore values. A semaphore is a positive integer. Since many applications require the use of  
3498 more than one semaphore, XSI-conformant systems have the ability to create sets or arrays of  
3499 semaphores.

3500 Calls to support semaphores include:

3501 *semctl()*, *semget()*, *semop()*

3502 Semaphore sets are created by using the *semget()* function.

3503 The message type of IPC allows process to communicate through the exchange of data stored in  
3504 buffers. This data is transmitted between processes in discrete portions known as messages.

3505 Calls to support message queues include:

3506 *msgctl()*, *msgget()*, *msgrcv()*, *msgsnd()*

3507 The share memory type of IPC allows two or more processes to share memory and consequently  
3508 the data contained therein. This is done by allowing processes to set up access to a common  
3509 memory address space. This sharing of memory provides a fast means of exchange of data  
3510 between processes.

3511 Calls to support shared memory include:

3512 *shmctl()*, *shmdt()*, *shmget()*

3513 The *ftok()* interface is also provided.

3514 **B.2.8 Realtime**3515 **Advisory Information**

3516 POSIX.1b contains an Informative Annex with proposed interfaces for “real-time files”. These  
 3517 interfaces could determine groups of the exact parameters required to do “direct I/O” or  
 3518 “extents”. These interfaces were objected to by a significant portion of the balloting group as too  
 3519 complex. A portable application had little chance of correctly navigating the large parameter  
 3520 space to match its desires to the system. In addition, they only applied to a new type of file  
 3521 (realtime files) and they told the implementation exactly what to do as opposed to advising the  
 3522 implementation on application behavior and letting it optimize for the system the (portable)  
 3523 application was running on. For example, it was not clear how a system that had a disk array  
 3524 should set its parameters.

3525 There seemed to be several overall goals:

- 3526 • Optimizing sequential access
- 3527 • Optimizing caching behavior
- 3528 • Optimizing I/O data transfer
- 3529 • Preallocation

3530 The advisory interfaces, *posix\_fadvise()* and *posix\_madvise()*, satisfy the first two goals. The  
 3531 POSIX\_FADV\_SEQUENTIAL and POSIX\_MADV\_SEQUENTIAL advice tells the  
 3532 implementation to expect serial access. Typically the system will prefetch the next several serial  
 3533 accesses in order to overlap I/O. It may also free previously accessed serial data if memory is  
 3534 tight. If the application is not doing serial access it can use POSIX\_FADV\_WILLNEED and  
 3535 POSIX\_MADV\_WILLNEED to accomplish I/O overlap, as required. When the application  
 3536 advises POSIX\_FADV\_RANDOM or POSIX\_MADV\_RANDOM behavior, the implementation  
 3537 usually tries to fetch a minimum amount of data with each request and it does not expect much  
 3538 locality. POSIX\_FADV\_DONTNEED and POSIX\_MADV\_DONTNEED allow the system to free  
 3539 up caching resources as the data will not be required in the near future.

3540 POSIX\_FADV\_NOREUSE tells the system that caching the specified data is not optimal. For file  
 3541 I/O, the transfer should go directly to the user buffer instead of being cached internally by the  
 3542 implementation. To portably perform direct disk I/O on all systems, the application must  
 3543 perform its I/O transfers according to the following rules:

- 3544 1. The user buffer should be aligned according to the {POSIX\_REC\_XFER\_ALIGN} *pathconf()*  
 3545 variable.
- 3546 2. The number of bytes transferred in an I/O operation should be a multiple of the  
 3547 {POSIX\_ALLOC\_SIZE\_MIN} *pathconf()* variable.
- 3548 3. The offset into the file at the start of an I/O operation should be a multiple of the  
 3549 {POSIX\_ALLOC\_SIZE\_MIN} *pathconf()* variable.
- 3550 4. The application should ensure that all threads which open a given file specify  
 3551 POSIX\_FADV\_NOREUSE to be sure that there is no unexpected interaction between  
 3552 threads using buffered I/O and threads using direct I/O to the same file.

3553 In some cases, a user buffer must be properly aligned in order to be transferred directly to/from  
 3554 the device. The {POSIX\_REC\_XFER\_ALIGN} *pathconf()* variable tells the application the proper  
 3555 alignment.

3556 The preallocation goal is met by the space control function, *posix\_fallocate()*. The application can  
 3557 use *posix\_fallocate()* to guarantee no [ENOSPC] errors and to improve performance by prepaying

3558 any overhead required for block allocation.

3559 Implementations may use information conveyed by a previous *posix\_fadvise()* call to influence  
3560 the manner in which allocation is performed. For example, if an application did the following  
3561 calls:

```
3562     fd = open("file");
3563     posix_fadvise(fd, offset, len, POSIX_FADV_SEQUENTIAL);
3564     posix_fallocate(fd, len, size);
```

3565 an implementation might allocate the file contiguously on disk.

3566 Finally, the *pathconf()* variables {*POSIX\_REC\_MIN\_XFER\_SIZE*},  
3567 {*POSIX\_REC\_MAX\_XFER\_SIZE*}, and {*POSIX\_REC\_INCR\_XFER\_SIZE*} tell the application a  
3568 range of transfer sizes that are recommended for best I/O performance.

3569 Where bounded response time is required, the vendor can supply the appropriate settings of the  
3570 advisories to achieve a guaranteed performance level.

3571 The interfaces meet the goals while allowing applications using regular files to take advantage of  
3572 performance optimizations. The interfaces tell the implementation expected application  
3573 behavior which the implementation can use to optimize performance on a particular system  
3574 with a particular dynamic load.

3575 The *posix\_memalign()* function was added to allow for the allocation of specifically aligned  
3576 buffers; for example, for {*POSIX\_REC\_XFER\_ALIGN*}.

3577 The working group also considered the alternative of adding a function which would return an  
3578 aligned pointer to memory within a user supplied buffer. This was not considered to be the best  
3579 method, because it potentially wastes large amounts of memory when buffers need to be aligned  
3580 on large alignment boundaries.

### 3581 **Message Passing**

3582 This section provides the rationale for the definition of the message passing interface in  
3583 IEEE Std. 1003.1-200x. This is presented in terms of the objectives, models, and requirements  
3584 imposed upon this interface.

#### 3585 • Objectives

3586 Many applications, including both realtime and database applications, require a means of  
3587 passing arbitrary amounts of data between cooperating processes comprising the overall  
3588 application on one or more processors. Many conventional interfaces for interprocess  
3589 communication are insufficient for realtime applications in that efficient and deterministic  
3590 data passing methods cannot be implemented. This has prompted the definition of message  
3591 passing interfaces providing these facilities:

- 3592 — Open a message queue.
- 3593 — Send a message to a message queue.
- 3594 — Receive a message from a queue, either synchronously or asynchronously.
- 3595 — Alter message queue attributes for flow and resource control.

3596 It is assumed that an application may consist of multiple cooperating processes and that  
3597 these processes may wish to communicate and coordinate their activities. The message  
3598 passing facility described in IEEE Std. 1003.1-200x allows processes to communicate through  
3599 system-wide queues. These message queues are accessed through names that may be path  
3600 names. A message queue can be opened for use by multiple sending and/or multiple

- 3601 receiving processes.
- 3602 • Background on Embedded Applications
- 3603 Interprocess communication utilizing message passing is a key facility for the construction of  
3604 deterministic, high-performance realtime applications. The facility is present in all realtime  
3605 systems and is the framework upon which the application is constructed. The performance of  
3606 the facility is usually a direct indication of the performance of the resulting application.
- 3607 Realtime applications, especially for embedded systems, are typically designed around the  
3608 performance constraints imposed by the message passing mechanisms. Applications for  
3609 embedded systems are typically very tightly constrained. Application writers expect to  
3610 design and control the entire system. In order to minimize system costs, the writer will  
3611 attempt to use all resources to their utmost and minimize the requirement to add additional  
3612 memory or processors.
- 3613 The embedded applications usually share address spaces and only a simple message passing  
3614 mechanism is required. The application can readily access common data incurring only  
3615 mutual-exclusion overheads. The models desired are the simplest possible with the  
3616 application building higher-level facilities only when needed.
- 3617 • Requirements
- 3618 The following requirements determined the features of the message passing facilities defined  
3619 in IEEE Std. 1003.1-200x:
- 3620 — Naming of Message Queues
- 3621 The mechanism for gaining access to a message queue is a path name evaluated in a  
3622 context that is allowed to be a file system name space, or it can be independent of any file  
3623 system. This is a specific attempt to allow implementations based on either method in  
3624 order to address both embedded systems and to also allow implementation in larger  
3625 systems.
- 3626 The interface of *mq\_open()* is defined to allow but not require the access control and name  
3627 conflicts resulting from utilizing a file system for name resolution. All required behavior  
3628 is specified for the access control case. Yet a conforming implementation, such as an  
3629 embedded system kernel, may define that there are no distinctions between users and  
3630 may define that all process have all access privileges.
- 3631 — Embedded System Naming
- 3632 Embedded systems need to be able to utilize independent name spaces for accessing the  
3633 various system objects. They typically do not have a file system, precluding its utilization  
3634 as a common name resolution mechanism. The modularity of an embedded system limits  
3635 the connections between separate mechanisms that can be allowed.
- 3636 Embedded systems typically do not have any access protection. Since the system does not  
3637 support the mixing of applications from different areas, and usually does not even have  
3638 the concept of an authorization entity, access control is not useful.
- 3639 — Large System Naming
- 3640 On systems with more functionality, the name resolution must support the ability to use  
3641 the file system as the name resolution mechanism/object storage medium and to have  
3642 control over access to the objects. Utilizing the path name space can result in further  
3643 errors when the names conflict with other objects.
- 3644 — Fixed Size of Messages



3645 The interfaces impose a fixed upper bound on the size of messages that can be sent to a  
3646 specific message queue. The size is set on an individual queue basis and cannot be  
3647 changed dynamically.

3648 The purpose of the fixed size is to increase the ability of the system to optimize the  
3649 implementation of *mq\_send()* and *mq\_receive()*. With fixed sizes of messages and fixed  
3650 numbers of messages, specific message blocks can be pre-allocated. This eliminates a  
3651 significant amount of checking for errors and boundary conditions. Additionally, an  
3652 implementation can optimize data copying to maximize performance. Finally, with a  
3653 restricted range of message sizes, an implementation is better able to provide  
3654 deterministic operations.

#### 3655 — Prioritization of Messages

3656 Message prioritization allows the application to determine the order in which messages  
3657 are received. Prioritization of messages is a key facility that is provided by most realtime  
3658 kernels and is heavily utilized by the applications. The major purpose of having priorities  
3659 in message queues is to avoid priority inversions in the message system, where a high-  
3660 priority message is delayed behind one or more lower-priority messages. It has been  
3661 observed that a significant problem with Ada rendezvous is that it queues tasks in strict  
3662 FIFO order, ignoring priorities. This allows the applications to be designed so that they do  
3663 not need to be interrupted in order to change the flow of control when exceptional  
3664 conditions occur. The prioritization does add additional overhead to the message  
3665 operations in those cases it is actually used but a clever implementation can optimize for  
3666 the FIFO case to make that more efficient.

#### 3667 — Asynchronous Notification

3668 The interface supports the ability to have a task asynchronously notified of the  
3669 availability of a message on the queue. The purpose of this facility is to allow the task to  
3670 perform other functions and yet still be notified that a message has become available on  
3671 the queue.

3672 To understand the requirement for this function, it is useful to understand two models of  
3673 application design: a single task performing multiple functions and multiple tasks  
3674 performing a single function. Each of these models has advantages.

3675 Asynchronous notification is required to build the model of a single task performing  
3676 multiple operations. This model typically results from either the expectation that  
3677 interruption is less expensive than utilizing a separate task or from the growth of the  
3678 application to include additional functions.

### 3679 **Semaphores**

3680 Semaphores are a high-performance process synchronization mechanism. Semaphores are  
3681 named by null-terminated strings of characters.

3682 A semaphore is created using the *sem\_init()* function or the *sem\_open()* function with the  
3683 *O\_CREAT* flag set in *oflag*.

3684 To use a semaphore, a process has to first initialize the semaphore or inherit an open descriptor  
3685 for the semaphore via *fork()*.

3686 A semaphore preserves its state when the last reference is closed. For example, if a semaphore  
3687 has a value of 13 when the last reference is closed, it will have a value of 13 when it is next  
3688 opened.

3689 When a semaphore is created, an initial state for the semaphore has to be provided. This value is  
3690 a non-negative integer. Negative values are not possible since they indicate the presence of  
3691 blocked processes. The persistence of any of these objects across a system crash or a system  
3692 reboot is undefined. Conforming applications shall not depend on any sort of persistence across  
3693 a system reboot or a system crash.

3694 • Models and Requirements

3695 A realtime system requires synchronization and communication between the processes  
3696 comprising the overall application. An efficient and reliable synchronization mechanism has  
3697 to be provided in a realtime system that will allow more than one schedulable process  
3698 mutually-exclusive access to the same resource. This synchronization mechanism has to  
3699 allow for the optimal implementation of synchronization or systems implementors will  
3700 define other, more cost-effective methods.

3701 At issue are the methods whereby multiple processes (tasks) can be designed and  
3702 implemented to work together in order to perform a single function. This requires  
3703 interprocess communication and synchronization. A semaphore mechanism is the lowest  
3704 level of synchronization that can be provided by an operating system.

3705 A semaphore is defined as an object that has an integral value and a set of blocked processes  
3706 associated with it. If the value is positive or zero, then the set of blocked processes is empty;  
3707 otherwise, the size of the set is equal to the absolute value of the semaphore value. The value  
3708 of the semaphore can be incremented or decremented by any process with access to the  
3709 semaphore and must be done as an indivisible operation. When a semaphore value is less  
3710 than or equal to zero, any process that attempts to lock it again will block or be informed that  
3711 it is not possible to perform the operation.

3712 A semaphore may be used to guard access to any resource accessible by more than one  
3713 schedulable task in the system. It is a global entity and not associated with any particular  
3714 process. As such, a method of obtaining access to the semaphore has to be provided by the  
3715 operating system. A process that wants access to a critical resource (section) has to wait on  
3716 the semaphore that guards that resource. When the semaphore is locked on behalf of a  
3717 process, it knows that it can utilize the resource without interference by any other  
3718 cooperating process in the system. When the process finishes its operation on the resource,  
3719 leaving it in a well-defined state, it posts the semaphore, indicating that some other process  
3720 may now obtain the resource associated with that semaphore.

3721 In this section, mutexes and condition variables are specified as the synchronization  
3722 mechanisms between threads.

3723 These primitives are typically used for synchronizing threads that share memory in a single  
3724 process. However, this section provides an option allowing the use of these synchronization  
3725 interfaces and objects between processes that share memory, regardless of the method for  
3726 sharing memory.

3727 Much experience with semaphores shows that there are two distinct uses of synchronization:  
3728 locking, which is typically of short duration; and waiting, which is typically of long or  
3729 unbounded duration. These distinct usages map directly onto mutexes and condition  
3730 variables, respectively.

3731 Semaphores are provided in IEEE Std. 1003.1-200x primarily to provide a means of  
3732 synchronization for processes; these processes may or may not share memory. Mutexes and  
3733 condition variables are specified as synchronization mechanisms between threads; these  
3734 threads always share (some) memory. Both are synchronization paradigms that have been in  
3735 widespread use for a number of years. Each set of primitives is particularly well matched to  
3736 certain problems.

3737 With respect to binary semaphores, experience has shown that condition variables and  
 3738 mutexes are easier to use for many synchronization problems than binary semaphores. The  
 3739 primary reason for this is the explicit appearance of a Boolean predicate that specifies when  
 3740 the condition wait is satisfied. This Boolean predicate terminates a loop, including the call to  
 3741 *pthread\_cond\_wait()*. As a result, extra wakeups are benign since the predicate governs  
 3742 whether the thread will actually proceed past the condition wait. With stateful primitives,  
 3743 such as binary semaphores, the wakeup in itself typically means that the wait is satisfied. The  
 3744 burden of ensuring correctness for such waits is thus placed on *all* signalers of the semaphore  
 3745 rather than on an *explicitly coded* Boolean predicate located at the condition wait. Experience  
 3746 has shown that the latter creates a major improvement in safety and ease-of-use.

3747 Counting semaphores are well matched to dealing with producer/consumer problems,  
 3748 including those that might exist between threads of different processes, or between a signal  
 3749 handler and a thread. In the former case, there may be little or no memory shared by the  
 3750 processes; in the latter case, one is not communicating between co-equal threads, but  
 3751 between a thread and an interruptlike entity. It is for these reasons that IEEE Std. 1003.1-200x  
 3752 allows semaphores to be used by threads.

3753 Mutexes and condition variables have been effectively used with and without priority  
 3754 inheritance, priority ceiling, and other attributes to synchronize threads that share memory.  
 3755 The efficiency of their implementation is comparable to or better than that of other  
 3756 synchronization primitives that are sometimes harder to use (for example, binary  
 3757 semaphores). Furthermore, there is at least one known implementation of Ada tasking that  
 3758 uses these primitives. Mutexes and condition variables together constitute an appropriate,  
 3759 sufficient, and complete set of interthread synchronization primitives.

3760 Efficient multi-threaded applications require high-performance synchronization primitives.  
 3761 Considerations of efficiency and generality require a small set of primitives upon which more  
 3762 sophisticated synchronization functions can be built.

#### 3763 • Standardization Issues

3764 It is possible to implement very high-performance semaphores using test-and-set  
 3765 instructions on shared memory locations. The library routines that implement such a high-  
 3766 performance interface has to properly ensure that a *sem\_wait()* or *sem\_trywait()* operation  
 3767 that cannot be performed will issue a blocking semaphore system call or properly report the  
 3768 condition to the application. The same interface to the application program would be  
 3769 provided by a high-performance implementation.

### 3770 B.2.8.1 Realtime Signals

#### 3771 **Realtime Signals Extension**

3772 This portion of the rationale presents models, requirements, and standardization issues relevant  
 3773 to the Realtime Signals Extension. This extension provides the capability required to support  
 3774 reliable, deterministic, asynchronous notification of events. While a new mechanism,  
 3775 unencumbered by the historical usage and semantics of POSIX.1 signals, might allow for a more  
 3776 efficient implementation, the application requirements for event notification can be met with a  
 3777 small number of extensions to signals. Therefore, a minimal set of extensions to signals to  
 3778 support the application requirements is specified.

3779 The realtime signal extensions specified in this section are used by other realtime functions  
 3780 requiring asynchronous notification:

#### 3781 • Models

3782 The model supported is one of multiple cooperating processes, each of which handles  
3783 multiple asynchronous external events. Events represent occurrences that are generated as  
3784 the result of some activity in the system. Examples of occurrences that can constitute an  
3785 event include:

3786 — Completion of an asynchronous I/O request

3787 — Expiration of a POSIX.1b timer

3788 — Arrival of an interprocess message

3789 — Generation of a user-defined event

3790 Processing of these events may occur synchronously via polling for event notifications or  
3791 asynchronously via a software interrupt mechanism. Existing practice for this model is well  
3792 established for traditional proprietary realtime operating systems, realtime executives, and  
3793 realtime extended POSIX-like systems.

3794 A contrasting model is that of “cooperating sequential processes” where each process  
3795 handles a single priority of events via polling. Each process blocks while waiting for events,  
3796 and each process depends on the preemptive, priority-based process scheduling mechanism  
3797 to arbitrate between events of different priority that need to be processed concurrently.  
3798 Existing practice for this model is also well established for small realtime executives that  
3799 typically execute in an unprotected physical address space, but it is just emerging in the  
3800 context of a fuller function operating system with multiple virtual address spaces.

3801 It could be argued that the cooperating sequential process model, and the facilities supported  
3802 by the POSIX Threads Extension obviate a software interrupt model. But, even with the  
3803 cooperating sequential process model, the need has been recognized for a software interrupt  
3804 model to handle exceptional conditions and process aborting, so the mechanism must be  
3805 supported in any case. Furthermore, it is not the purview of IEEE Std. 1003.1-200x to attempt  
3806 to convince realtime practitioners that their current application models based on software  
3807 interrupts are “broken” and should be replaced by the cooperating sequential process model.  
3808 Rather, it is the charter of IEEE Std. 1003.1-200x to provide standard extensions to  
3809 mechanisms that support existing realtime practice.

3810 • Requirements

3811 This section discusses the following realtime application requirements for asynchronous  
3812 event notification:

3813 — Reliable delivery of asynchronous event notification

3814 The events notification mechanism shall guarantee delivery of an event notification.  
3815 Asynchronous operations (such as asynchronous I/O and timers) that complete  
3816 significantly after they are invoked have to guarantee that delivery of the event  
3817 notification can occur at the time of completion.

3818 — Prioritized handling of asynchronous event notifications

3819 The events notification mechanism shall support the assigning of a user function as an  
3820 event notification handler. Furthermore, the mechanism shall support the preemption of  
3821 an event handler function by a higher priority event notification and shall support the  
3822 selection of the highest priority pending event notification when multiple notifications (of  
3823 different priority) are pending simultaneously.

3824 The model here is based on hardware interrupts. Asynchronous event handling allows  
3825 the application to ensure that time-critical events are immediately processed when  
3826 delivered, without the indeterminism of being at a random location within a polling loop.

- 3827 Use of handler priority allows the specification of how handlers are interrupted by other  
3828 higher priority handlers.
- 3829 — Differentiation between multiple occurrences of event notifications of the same type
- 3830 The events notification mechanism shall pass an application-defined value to the event  
3831 handler function. This value can be used for a variety of purposes, such as enabling the  
3832 application to identify which of several possible events of the same type (for example,  
3833 timer expirations) has occurred.
- 3834 — Polled reception of asynchronous event notifications
- 3835 The events notification mechanism shall support blocking and non-blocking polls for  
3836 asynchronous event notification.
- 3837 The polled mode of operation is often preferred over the interrupt mode by those  
3838 practitioners accustomed to this model. Providing support for this model facilitates the  
3839 porting of applications based on this model to POSIX.1b conforming systems.
- 3840 — Deterministic response to asynchronous event notifications
- 3841 The events notification mechanism shall not preclude implementations that provide  
3842 deterministic event dispatch latency and shall minimize the number of system calls  
3843 needed to use the event facilities during realtime processing.
- 3844 • Rationale for Extension
- 3845 POSIX.1 signals have many of the characteristics necessary to support the asynchronous  
3846 handling of event notifications, and the Realtime Signals Extension addresses the following  
3847 deficiencies in the POSIX.1 signal mechanism:
- 3848 — Signals do not support reliable delivery of event notification. Subsequent occurrences of  
3849 a pending signal are not guaranteed to be delivered.
- 3850 — Signals do not support prioritized delivery of event notifications. The order of signal  
3851 delivery when multiple unblocked signals are pending is undefined.
- 3852 — Signals do not support the differentiation between multiple signals of the same type.

### 3853 *B.2.8.2 Asynchronous I/O*

3854 Many applications need to interact with the I/O subsystem in an asynchronous manner. The  
3855 asynchronous I/O mechanism provides the ability to overlap application processing and I/O  
3856 operations initiated by the application. The asynchronous I/O mechanism allows a single  
3857 process to perform I/O simultaneously to a single file multiple times or to multiple files  
3858 multiple times.

#### 3859 **Overview**

3860 Asynchronous I/O operations proceed in logical parallel with the processing done by the  
3861 application after the asynchronous I/O has been initiated. Other than this difference,  
3862 asynchronous I/O behaves similarly to normal I/O using *read()*, *write()*, *lseek()*, and *fsync()*.  
3863 The effect of issuing an asynchronous I/O request is as if a separate thread of execution were to  
3864 perform atomically the implied *lseek()* operation, if any, and then the requested I/O operation  
3865 (either *read()*, *write()*, or *fsync()*). There is no seek implied with a call to *aio\_fsync()*. Concurrent  
3866 asynchronous operations and synchronous operations applied to the same file update the file as  
3867 if the I/O operations had proceeded serially.

3868 When asynchronous I/O completes, a signal can be delivered to the application to indicate the  
3869 completion of the I/O. This signal can be used to indicate that buffers and control blocks used

3870 for asynchronous I/O can be reused. Signal delivery is not required for an asynchronous  
 3871 operation and may be turned off on a per-operation basis by the application. Signals may also be  
 3872 synchronously polled using *aio\_suspend()*, *sigtimedwait()*, or *sigwaitinfo()*.

3873 Normal I/O has a return value and an error status associated with it. Asynchronous I/O returns  
 3874 a value and an error status when the operation is first submitted, but that only relates to whether  
 3875 the operation was successfully queued up for servicing. The I/O operation itself also has a  
 3876 return status and an error value. To allow the application to retrieve the return status and the  
 3877 error value, functions are provided that, given the address of an asynchronous I/O control  
 3878 block, yield the return and error status associated with the operation. Until an asynchronous I/O  
 3879 operation is done, its error status shall be [EINPROGRESS]. Thus, an application can poll for  
 3880 completion of an asynchronous I/O operation by waiting for the error status to become equal to  
 3881 a value other than [EINPROGRESS]. The return status of an asynchronous I/O operation is  
 3882 undefined so long as the error status is equal to [EINPROGRESS].

3883 Storage for asynchronous operation return and error status may be limited. Submission of  
 3884 asynchronous I/O operations may fail if this storage is exceeded. When an application retrieves  
 3885 the return status of a given asynchronous operation, therefore, any system-maintained storage  
 3886 used for this status and the error status may be reclaimed for use by other asynchronous  
 3887 operations.

3888 Asynchronous I/O can be performed on file descriptors that have been enabled for POSIX.1b  
 3889 synchronized I/O. In this case, the I/O operation still occurs asynchronously, as defined herein;  
 3890 however, the asynchronous operation I/O in this case is not completed until the I/O has reached  
 3891 either the state of synchronized I/O data integrity completion or synchronized I/O file integrity  
 3892 completion, depending on the sort of synchronized I/O that is enabled on the file descriptor.

### 3893 **Models**

3894 Three models illustrate the use of asynchronous I/O: a journalization model, a data acquisition  
 3895 model, and a model of the use of asynchronous I/O in supercomputing applications.

- 3896 • **Journalization Model**

3897 Many realtime applications perform low-priority journalizing functions. Journalizing  
 3898 requires that logging records be queued for output without blocking the initiating process.

- 3899 • **Data Acquisition Model**

3900 A data acquisition process may also serve as a model. The process has two or more channels  
 3901 delivering intermittent data that must be read within a certain time. The process issues one  
 3902 asynchronous read on each channel. When one of the channels needs data collection, the  
 3903 process reads the data and posts it through an asynchronous write to secondary memory for  
 3904 future processing.

- 3905 • **Supercomputing Model**

3906 The supercomputing community has used asynchronous I/O much like that specified herein  
 3907 for many years. This community requires the ability to perform multiple I/O operations to  
 3908 multiple devices with a minimal number of entries to “the system”; each entry to “the  
 3909 system” provokes a major delay in operations when compared to the normal progress made  
 3910 by the application. This existing practice motivated the use of combined *lseek()* and *read()* or  
 3911 *write()* calls, as well as the *lio\_listio()* call. Another common practice is to disable signal  
 3912 notification for I/O completion, and simply poll for I/O completion at some interval by  
 3913 which the I/O should be completed. Likewise, interfaces like *aio\_cancel()* have been in  
 3914 successful commercial use for many years. Note also that an underlying implementation of  
 3915 asynchronous I/O will require the ability, at least internally, to cancel outstanding

3916 asynchronous I/O, at least when the process exits. (Consider an asynchronous read from a  
3917 terminal, when the process intends to exit immediately.)

### 3918 **Requirements**

3919 Asynchronous input and output for realtime implementations have these requirements:

- 3920 • The ability to queue multiple asynchronous read and write operations to a single open  
3921 instance. Both sequential and random access should be supported.
- 3922 • The ability to queue asynchronous read and write operations to multiple open instances.
- 3923 • The ability to obtain completion status information by polling and/or asynchronous event  
3924 notification.
- 3925 • Asynchronous event notification on asynchronous I/O completion is optional.
- 3926 • It has to be possible for the application to associate the event with the *aiocbp* for the operation  
3927 that generated the event.
- 3928 • The ability to cancel queued requests.
- 3929 • The ability to wait upon asynchronous I/O completion in conjunction with other types of  
3930 events.
- 3931 • The ability to accept an *aio\_read()* and an *aio\_cancel()* for a device that accepts a *read()*, and  
3932 the ability to accept an *aio\_write()* and an *aio\_cancel()* for a device that accepts a *write()*. This  
3933 does not imply that the operation is asynchronous.

### 3934 **Standardization Issues**

3935 The following issues are addressed by the standardization of asynchronous I/O:

- 3936 • Rationale for New Interface

3937 Non-blocking I/O does not satisfy the needs of either realtime or high-performance  
3938 computing models; these models require that a process overlap program execution and I/O  
3939 processing. Realtime applications will often make use of direct I/O to or from the address  
3940 space of the process, or require synchronized (unbuffered) I/O; they also require the ability  
3941 to overlap this I/O with other computation. In addition, asynchronous I/O allows an  
3942 application to keep a device busy at all times, possibly achieving greater throughput.  
3943 Supercomputing and database architectures will often have specialized hardware that can  
3944 provide true asynchrony underlying the logical asynchrony provided by this interface. In  
3945 addition, asynchronous I/O should be supported by all types of files and devices in the same  
3946 manner.

- 3947 • Effect of Buffering

3948 If asynchronous I/O is performed on a file that is buffered prior to being actually written to  
3949 the device, it is possible that asynchronous I/O will offer no performance advantage over  
3950 normal I/O; the cycles *stolen* to perform the asynchronous I/O will be taken away from the  
3951 running process and the I/O will occur at interrupt time. This potential lack of gain in  
3952 performance in no way obviates the need for asynchronous I/O by realtime applications,  
3953 which very often will use specialized hardware support; multiple processors; and/or  
3954 unbuffered, synchronized I/O.

3955 **B.2.8.3** *Memory Management*

3956 All memory management and shared memory definitions are located in the `<sys/mman.h>`  
3957 header. This is for alignment with historical practice.

3958 **Memory Locking Functions**

3959 This portion of the rationale presents models, requirements, and standardization issues relevant  
3960 to process memory locking.

## 3961 • Models

3962 Realtime systems that conform to IEEE Std. 1003.1-200x are expected (and desired) to be  
3963 supported on systems with demand-paged virtual memory management, non-paged  
3964 swapping memory management, and physical memory systems with no memory  
3965 management hardware. The general case, however, is the demand-paged, virtual memory  
3966 system with each POSIX process running in a virtual address space. Note that this includes  
3967 architectures where each process resides in its own virtual address space and architectures  
3968 where the address space of each process is only a portion of a larger global virtual address  
3969 space.

3970 The concept of memory locking is introduced to eliminate the indeterminacy introduced by  
3971 paging and swapping, and to support an upper bound on the time required to access the  
3972 memory mapped into the address space of a process. Ideally, this upper bound will be the  
3973 same as the time required for the processor to access “main memory”, including any address  
3974 translation and cache miss overheads. But some implementations—primarily on  
3975 mainframes—will not actually force locked pages to be loaded and held resident in main  
3976 memory. Rather, they will handle locked pages so that accesses to these pages will meet the  
3977 performance metrics for locked process memory in the implementation. Also, although it is  
3978 not, for example, the intention that this interface, as specified, be used to lock process  
3979 memory into “cache”, it is conceivable that an implementation could support a large static  
3980 RAM memory and define this as “main memory” and use a large[r] dynamic RAM as  
3981 “backing store”. These interfaces could then be interpreted as supporting the locking of  
3982 process memory into the static RAM. Support for multiple levels of backing store would  
3983 require extensions to these interfaces.

3984 Implementations may also use memory locking to guarantee a fixed translation between  
3985 virtual and physical addresses where such is beneficial to improving determinacy for  
3986 direct-to/from-process input/output. IEEE Std. 1003.1-200x does not guarantee to the  
3987 application that the virtual-to-physical address translations, if such exist, are fixed, because  
3988 such behavior would not be implementable on all architectures on which implementations of  
3989 IEEE Std. 1003.1-200x are expected. But IEEE Std. 1003.1-200x does mandate that an  
3990 implementation define, for the benefit of potential users, whether or not locking guarantees  
3991 fixed translations.

3992 Memory locking is defined with respect to the address space of a process. Only the pages  
3993 mapped into the address space of a process may be locked by the process, and when the  
3994 pages are no longer mapped into the address space—for whatever reason—the locks  
3995 established with respect to that address space are removed. Shared memory areas warrant  
3996 special mention, as they may be mapped into more than one address space or mapped more  
3997 than once into the address space of a process; locks may be established on pages within these  
3998 areas with respect to several of these mappings. In such a case, the lock state of the  
3999 underlying physical pages is the logical OR of the lock state with respect to each of the  
4000 mappings. Only when all such locks have been removed are the shared pages considered  
4001 unlocked.



4002 In recognition of the page granularity of Memory Management Units (MMU), and in order to  
4003 support locking of ranges of address space, memory locking is defined in terms of “page”  
4004 granularity. That is, for the interfaces that support an address and size specification for the  
4005 region to be locked, the address must be on a page boundary, and all pages mapped by the  
4006 specified range are locked, if valid. This means that the length is implicitly rounded up to a  
4007 multiple of the page size. The page size is implementation-defined and is available to  
4008 applications as a compile time symbolic constant or at runtime via *sysconf()*.

4009 A “real memory” POSIX.1b implementation that has no MMU could elect not to support  
4010 these interfaces, returning [ENOSYS]. But an application could easily interpret this as  
4011 meaning that the implementation would unconditionally page or swap the application when  
4012 such is not the case. It is the intention of IEEE Std. 1003.1-200x that such a system could  
4013 define these interfaces as “NO-OPs”, returning success without actually performing any  
4014 function except for mandated argument checking.

4015 • Requirements

4016 For realtime applications, memory locking is generally considered to be required as part of  
4017 application initialization. This locking is performed after an application has been loaded (that  
4018 is, *exec'd*) and the program remains locked for its entire lifetime. But to support applications  
4019 that undergo major mode changes where, in one mode, locking is required, but in another it  
4020 is not, the specified interfaces allow repeated locking and unlocking of memory within the  
4021 lifetime of a process.

4022 When a realtime application locks its address space, it should not be necessary for the  
4023 application to then “touch” all of the pages in the address space to guarantee that they are  
4024 resident or else suffer potential paging delays the first time the page is referenced. Thus,  
4025 IEEE Std. 1003.1-200x requires that the pages locked by the specified interfaces be resident  
4026 when the locking functions return successfully.

4027 Many architectures support system-managed stacks that grow automatically when the  
4028 current extent of the stack is exceeded. A realtime application has a requirement to be able to  
4029 “preallocate” sufficient stack space and lock it down so that it will not suffer page faults to  
4030 grow the stack during critical realtime operation. There was no consensus on a portable way  
4031 to specify how much stack space is needed, so IEEE Std. 1003.1-200x supports no specific  
4032 interface for preallocating stack space. But an application can portably lock down a specific  
4033 amount of stack space by specifying *MCL\_FUTURE* in a call to *memlockall()* and then calling  
4034 a dummy function that declares an automatic array of the desired size.

4035 Memory locking for realtime applications is also generally considered to be an “all or  
4036 nothing” proposition. That is, the entire process, or none, is locked down. But, for  
4037 applications that have well-defined sections that need to be locked and others that do not,  
4038 IEEE Std. 1003.1-200x supports an optional set of interfaces to lock or unlock a range of  
4039 process addresses. Reasons for locking down a specific range include:

4040 — An asynchronous event handler function that must respond to external events in a  
4041 deterministic manner such that page faults cannot be tolerated

4042 — An input/output “buffer” area that is the target for direct-to-process I/O, and the  
4043 overhead of implicit locking and unlocking for each I/O call cannot be tolerated

4044 Finally, locking is generally viewed as an “application-wide” function. That is, the  
4045 application is globally aware of which regions are locked and which are not over time. This is  
4046 in contrast to a function that is used temporarily within a “third party” library routine whose  
4047 function is unknown to the application, and therefore must have no “side effects”. The  
4048 specified interfaces, therefore, do not support “lock stacking” or “lock nesting” within a  
4049 process. But, for pages that are shared between processes or mapped more than once into a

4050 process address space, “lock stacking” is essentially mandated by the requirement that  
4051 unlocking of pages that are mapped by more than one process or more than once by the same  
4052 process does not affect locks established on the other mappings.

4053 There was some support for “lock stacking” so that locking could be transparently used in  
4054 library functions or opaque modules. But the consensus was not to burden all  
4055 implementations with lock stacking (and reference counting), and an implementation option  
4056 was proposed. There were strong objections to the option because applications would have  
4057 to support both options in order to remain portable. The consensus was to eliminate lock  
4058 stacking altogether, primarily through overwhelming support for the System V  
4059 “m[un]lock[all]” interface on which IEEE Std. 1003.1-200x is now based.

4060 Locks are not inherited across *fork()*s because some systems implement *fork()* by creating  
4061 new address spaces for the child. In such an implementation, requiring locks to be inherited  
4062 would lead to new situations in which a fork would fail due to the inability of the system to  
4063 lock sufficient memory to lock both the parent and the child. The consensus was that there  
4064 was no benefit to such inheritance. Note that this does not mean that locks are removed  
4065 when, for instance, a thread is created in the same address space.

4066 Similarly, locks are not inherited across *exec* because some systems implement *exec* by  
4067 unmapping all of the pages in the address space (which, by definition, removes the locks on  
4068 these pages), and maps in pages of the *exec*'d image. In such an implementation, requiring  
4069 locks to be inherited would lead to new situations in which *exec* would fail. Reporting this  
4070 failure would be very cumbersome to detect in time to report to the calling process, and no  
4071 appropriate mechanism exists for informing the *exec*'d process of its status.

4072 It was determined that, if the newly loaded application required locking, it was the  
4073 responsibility of that application to establish the locks. This is also in keeping with the  
4074 general view that it is the responsibility of the application to be aware of all locks that are  
4075 established.

4076 There was one request to allow (not mandate) locks to be inherited across *fork()*, and a  
4077 request for a flag, MCL\_INHERIT, that would specify inheritance of memory locks across  
4078 *exec*s. Given the difficulties raised by this and the general lack of support for the feature in  
4079 IEEE Std. 1003.1-200x, it was not added. IEEE Std. 1003.1-200x does not preclude an  
4080 implementation from providing this feature for administrative purposes, such as a “run”  
4081 command that will lock down and execute specified program. Additionally, the rationale for  
4082 the objection equated *fork()* with creating a thread in the address space. IEEE Std. 1003.1-200x  
4083 does not mandate releasing locks when creating additional threads in an existing process.

#### 4084 • Standardization Issues

4085 One goal of IEEE Std. 1003.1-200x is to define a set of primitives that provide the necessary  
4086 functionality for realtime applications, with consideration for the needs of other application  
4087 domains where such were identified, which is based to the extent possible on existing  
4088 industry practice.

4089 The Memory Locking option is required by many realtime applications to tune performance.  
4090 Such a facility is accomplished by placing constraints on the virtual memory system to limit  
4091 paging of time of the process or of critical sections of the process. This facility should not be  
4092 used by most non-realtime applications.

4093 Optional features provided in IEEE Std. 1003.1-200x allow applications to lock selected  
4094 address ranges with the caveat that the process is responsible for being aware of the page  
4095 granularity of locking and the un-nested nature of the locks.

4096 **Mapped Files Functions**

4097 The Memory Mapped Files option provides a mechanism that allows a process to access files by  
4098 directly incorporating file data into its address space. Once a file is “mapped” into a process  
4099 address space, the data can be manipulated by instructions as memory. The use of mapped files  
4100 can significantly reduce I/O data movement since file data does not have to be copied into  
4101 process data buffers as in *read()* and *write()*. If more than one process maps a file, its contents  
4102 are shared among them. This provides a low overhead mechanism by which processes can  
4103 synchronize and communicate.

## 4104 • Historical Perspective

4105 Realtime applications have historically been implemented using a collection of cooperating  
4106 processes or tasks. In early systems, these processes ran on bare hardware (that is, without an  
4107 operating system) with no memory relocation or protection. The application paradigms that  
4108 arose from this environment involve the sharing of data between the processes.

4109 When realtime systems were implemented on top of vendor-supplied operating systems, the  
4110 paradigm or performance benefits of direct access to data by multiple processes was still  
4111 deemed necessary. As a result, operating systems that claim to support realtime applications  
4112 must support the shared memory paradigm.

4113 Additionally, a number of realtime systems provide the ability to map specific sections of the  
4114 physical address space into the address space of a process. This ability is required if an  
4115 application is to obtain direct access to memory locations that have specific properties (for  
4116 example, refresh buffers or display devices, dual ported memory locations, DMA target  
4117 locations). The use of this ability is common enough to warrant some degree of  
4118 standardization of its interface. This ability overlaps the general paradigm of shared  
4119 memory in that, in both instances, common global objects are made addressable by  
4120 individual processes or tasks.

4121 Finally, a number of systems also provide the ability to map process addresses to files. This  
4122 provides both a general means of sharing persistent objects, and using files in a manner that  
4123 optimizes memory and swapping space usage.

4124 Simple shared memory is clearly a special case of the more general file mapping capability.  
4125 In addition, there is relatively widespread agreement and implementation of the file  
4126 mapping interface. In these systems, many different types of objects can be mapped (for  
4127 example, files, memory, devices, and so on) using the same mapping interfaces. This  
4128 approach both minimizes interface proliferation and maximizes the generality of programs  
4129 using the mapping interfaces.

## 4130 • Memory Mapped Files Usage

4131 A memory object can be concurrently mapped into the address space of one or more  
4132 processes. The *mmap()* and *munmap()* functions allow a process to manipulate their address  
4133 space by mapping portions of memory objects into it and removing them from it. When  
4134 multiple processes map the same memory object, they can share access to the underlying  
4135 data. Implementations may restrict the size and alignment of mappings to be on *page*-size  
4136 boundaries. The page size, in bytes, is the value of the system-configurable variable  
4137 {PAGESIZE}, typically accessed by calling *sysconf()* with a *name* argument of  
4138 *\_SC\_PAGESIZE*. If an implementation has no restrictions on size or alignment, it may  
4139 specify a 1-byte page size.

4140 To map memory, a process first opens a memory object. The *ftruncate()* function can be used  
4141 to contract or extend the size of the memory object even when the object is currently  
4142 mapped. If the memory object is extended, the contents of the extended areas are zeros.

4143 After opening a memory object, the application maps the object into its address space using  
4144 the *mmap()* function call. Once a mapping has been established, it remains mapped until  
4145 unmapped with *munmap()*, even if the memory object is closed. The *mprotect()* function can  
4146 be used to change the memory protections initially established by *mmap()*.

4147 A *close()* of the file descriptor, while invalidating the file descriptor itself, does not unmap  
4148 any mappings established for the memory object. The address space, including all mapped  
4149 regions, is inherited on *fork()*. The entire address space is unmapped on process termination  
4150 or by successful calls to any of the *exec* family of functions.

4151 The *msync()* function is used to force mapped file data to permanent storage.

4152 • Effects on Other Functions

4153 When the Memory Mapped Files option is supported, the operation of the *open()*, *creat()*, and  
4154 *unlink()* functions are a natural result of using the file system name space to map the global  
4155 names for memory objects.

4156 The *ftruncate()* function can be used to set the length of a sharable memory object.

4157 The meaning of *stat()* fields other than the size and protection information is undefined on  
4158 implementations where memory objects are not implemented using regular files. When  
4159 regular files are used, the times reflect when the implementation updated the file image of  
4160 the data, not when a process updated the data in memory.

4161 The operations of *fdopen()*, *write()*, *read()*, and *lseek()* were made unspecified for objects  
4162 opened with *shm\_open()*, so that implementations that did not implement memory objects as  
4163 regular files would not have to support the operation of these functions on shared memory  
4164 objects.

4165 The behavior of memory objects with respect to *close()*, *dup()*, *dup2()*, *open()*, *close()*, *fork()*,  
4166 *\_exit()*, and the *exec* family of functions is the same as the behavior of the existing practice of  
4167 the *mmap()* function.

4168 A memory object can still be referenced after a close. That is, any mappings made to the file  
4169 are still in effect, and reads and writes that are made to those mappings are still valid and are  
4170 shared with other processes that have the same mapping. Likewise, the memory object can  
4171 still be used if any references remain after its name(s) have been deleted. Any references that  
4172 remain after a close must not appear to the application as file descriptors.

4173 This is existing practice for *mmap()* and *close()*. In addition, there are already mappings  
4174 present (text, data, stack) that do not have open file descriptors. The text mapping in  
4175 particular is considered a reference to the file containing the text. The desire was to treat all  
4176 mappings by the process uniformly. Also, many modern implementations use *mmap()* to  
4177 implement shared libraries, and it would not be desirable to keep file descriptors for each of  
4178 the many libraries an application can use. It was felt there were many other existing  
4179 programs that used this behavior to free a file descriptor, and thus IEEE Std. 1003.1-200x  
4180 could not forbid it and still claim to be using existing practice.

4181 For implementations that implement memory objects using memory only, memory objects  
4182 will retain the memory allocated to the file after the last close and will use that same memory  
4183 on the next open. Note that closing the memory object is not the same as deleting the name,  
4184 since the memory object is still defined in the memory object name space.

4185 The locks of *fcntl()* do not block any read or write operation, including read or write access to  
4186 shared memory or mapped files. In addition, implementations that only support shared  
4187 memory objects should not be required to implement record locks. The reference to *fcntl()* is  
4188 added to make this point explicitly. The other *fcntl()* commands are useful with shared

4189 memory objects.

4190 The size of pages that mapping hardware may be able to support may be a configurable  
4191 value, or it may change based on hardware implementations. The addition of the  
4192 `_SC_PAGESIZE` parameter to the `sysconf()` function is provided for determining the mapping  
4193 page size at runtime.

### 4194 Shared Memory Functions

4195 Implementations may support the Shared Memory Objects option without supporting a general  
4196 Memory Mapped Files option. Shared memory objects are named regions of storage that may be  
4197 independent of the file system and can be mapped into the address space of one or more  
4198 processes to allow them to share the associated memory.

#### 4199 • Requirements

4200 Shared memory is used to share data among several processes, each potentially running at  
4201 different priority levels, responding to different inputs, or performing separate tasks. Shared  
4202 memory is not just simply providing common access to data, it is providing the fastest  
4203 possible communication between the processes. With one memory write operation, a process  
4204 can pass information to as many processes as have the memory region mapped.

4205 As a result, shared memory provides a mechanism that can be used for all other interprocess  
4206 communications facilities. It may also be used by an application for implementing more  
4207 sophisticated mechanisms than semaphores and message queues.

4208 The need for a shared memory interface is obvious for virtual memory systems, where the  
4209 operating system is directly preventing processes from accessing each other's data. However,  
4210 in unprotected systems, such as those found in some embedded controllers, a shared  
4211 memory interface is needed to provide a portable mechanism to allocate a region of memory  
4212 to be shared and then to communicate the address of that region to other processes.

4213 This, then, provides the minimum functionality that a shared memory interface must have in  
4214 order to support realtime applications: to allocate and name an object to be mapped into  
4215 memory for potential sharing (`open()` or `shm_open()`), and to make the memory object  
4216 available within the address space of a process (`mmap()`). To complete the interface, a  
4217 mechanism to release the claim of a process on a shared memory object (`munmap()`) is also  
4218 needed, as well as a mechanism for deleting the name of a sharable object that was  
4219 previously created (`unlink()` or `shm_unlink()`).

4220 After a mapping has been established, an implementation should not have to provide  
4221 services to maintain that mapping. All memory writes into that area will appear immediately  
4222 in the memory mapping of that region by any other processes.

4223 Thus, requirements include:

4224 — Support creation of sharable memory objects and the mapping of these objects into the  
4225 address space of a process.

4226 — Sharable memory objects should be accessed by global names accessible from all  
4227 processes.

4228 — Support the mapping of specific sections of physical address space (such as a memory  
4229 mapped device) into the address space of a process. This should not be done by the  
4230 process specifying the actual address, but again by an implementation-defined global  
4231 name (such as a special device name) dedicated to this purpose.

4232 — Support the mapping of discrete portions of these memory objects.

- 4233 — Support for minimum hardware configurations that contain no physical media on which  
4234 to store shared memory contents permanently.
- 4235 — The ability to preallocate the entire shared memory region so that minimum hardware  
4236 configurations without virtual memory support can guarantee contiguous space.
- 4237 — The maximizing of performance by not requiring functionality that would require  
4238 implementation interaction above creating the shared memory area and returning the  
4239 mapping.
- 4240 Note that the above requirements do not preclude:
- 4241 — The sharable memory object from being implemented using actual files on an actual file  
4242 system.
- 4243 — The global name that is accessible from all processes being restricted to a file system area  
4244 that is dedicated to handling shared memory.
- 4245 — An implementation not providing implementation-defined global names for the purpose  
4246 of physical address mapping.
- 4247 • Shared Memory Objects Usage
- 4248 If the Shared Memory Objects option is supported, a shared memory object may be created,  
4249 or opened if it already exists, with the *shm\_open()* function. If the shared memory object is  
4250 created, it has a length of zero. The *truncate()* function can be used to set the size of the  
4251 shared memory object after creation. The *shm\_unlink()* function removes the name for a  
4252 shared memory object created by *shm\_open()*.
- 4253 • Shared Memory Overview
- 4254 The shared memory facility defined by IEEE Std. 1003.1-200x usually results in memory  
4255 locations being added to the address space of the process. The implementation returns the  
4256 address of the new space to the application by means of a pointer. This works well in  
4257 languages like C. However, in languages such as FORTRAN, it will not work because these  
4258 languages do not have pointer types. In the bindings for such a language, either a special  
4259 COMMON section will need to be defined (which is unlikely), or the binding will have to  
4260 allow existing structures to be mapped. The implementation will likely have to place  
4261 restrictions on the size and alignment of such structures or will have to map a suitable region  
4262 of the address space of the process into the memory object, and thus into other processes.  
4263 These are issues for that particular language binding. For IEEE Std. 1003.1-200x, however, the  
4264 practice will not be forbidden, merely undefined.
- 4265 Two potentially different name spaces are used for naming objects that may be mapped into  
4266 process address spaces. When the Memory Mapped Files option is supported, files may be  
4267 accessed via *open()*. When the Shared Memory Objects option is supported, sharable  
4268 memory objects that might not be files may be accessed via the *shm\_open()* function. These  
4269 options are not mutually-exclusive.
- 4270 Some systems supporting the Shared Memory Objects option may choose to implement the  
4271 shared memory object name space as part of the file system name space. There are several  
4272 reasons for this:
- 4273 — It allows applications to prevent name conflicts by use of the directory structure.
- 4274 — It uses an existing mechanism for accessing global objects and prevents the creation of a  
4275 new mechanism for naming global objects.
- 4276 In such implementations, memory objects can be implemented using regular files, if that is  
4277 what the implementation chooses. The *shm\_open()* function can be implemented as an *open()*

4278 call in a fixed directory followed by a call to *fcntl()* to set *FD\_CLOEXEC*. The *shm\_unlink()*  
 4279 function can be implemented as an *unlink()* call.

4280 On the other hand, it is also expected that small embedded systems that support the Shared  
 4281 Memory Objects option may wish to implement shared memory without having any file  
 4282 systems present. In this case, the implementations may choose to use a simple string valued  
 4283 name space for shared memory regions. The *shm\_open()* function permits either type of  
 4284 implementation.

4285 Some systems have hardware that supports protection of mapped data from certain classes  
 4286 of access and some do not. Systems that supply this functionality can support the Memory  
 4287 Protection option.

4288 Some implementations restrict size, alignment, and protections to be on *page-size*  
 4289 boundaries. If an implementation has no restrictions on size or alignment, it may specify a 1-  
 4290 byte page size. Applications on implementations that do support larger pages must be  
 4291 cognizant of the page size since this is the alignment and protection boundary.

4292 Simple embedded implementations may have a 1-byte page size and only support the Shared  
 4293 Memory Objects option. This provides simple shared memory between processes without  
 4294 requiring mapping hardware.

4295 IEEE Std. 1003.1-200x is silent about how implementations that chose to implement memory  
 4296 objects directly would treat them with standard utilities such as *ls*, because utilities are not  
 4297 within the charter of IEEE Std. 1003.1-200x.

4298 IEEE Std. 1003.1-200x specifically allows a memory object to remain referenced after a close  
 4299 because that is existing practice for the *mmap()* function.

#### 4300 **Typed Memory Functions**

4301 Implementations may support the Typed Memory Objects option without supporting either the  
 4302 Shared Memory option or the Memory Mapped Files option. Typed memory objects are pools of  
 4303 specialized storage, different from the main memory resource normally used by a processor to  
 4304 hold code and data, that can be mapped into the address space of one or more processes.

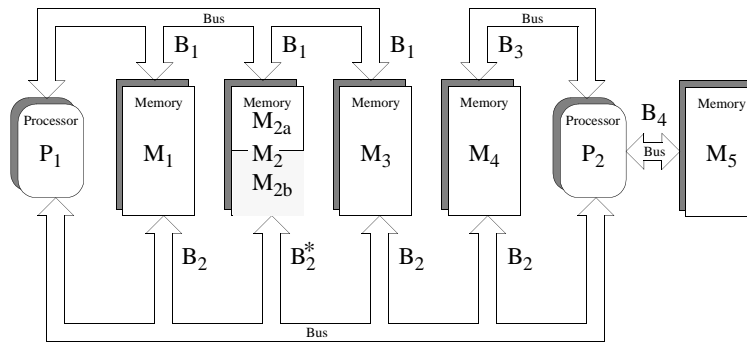
##### 4305 • Model

4306 Realtime systems conforming to one of the POSIX.13 realtime profiles are expected (and  
 4307 desired) to be supported on systems with more than one type or pool of memory (for  
 4308 example, SRAM, DRAM, ROM, EPROM, EEPROM), where each type or pool of memory may  
 4309 be accessible by one or more processors via one or more busses (ports). Memory mapped  
 4310 files, shared memory objects, and the language-specific storage allocation operators (*malloc()*  
 4311 for the ISO C standard, *new* for ANSI Ada) fail to provide application program interfaces  
 4312 versatile enough to allow applications to control their utilization of such diverse memory  
 4313 resources. The typed memory interfaces *posix\_typed\_mem\_open()*, *posix\_mem\_offset()*,  
 4314 *posix\_typed\_mem\_get\_info()*, *mmap()*, and *munmap()* defined herein support the model of  
 4315 typed memory described below.

4316 For purposes of this model, a system comprises several processors (for example, P1 and P2),  
 4317 several physical memory pools (for example, M1, M2, M2a, M2b, M3, M4, and M5), and  
 4318 several busses or “ports” (for example, B1, B2, B3, and B4) interconnecting the various  
 4319 processors and memory pools in some system-specific way. Notice that some memory pools  
 4320 may be contained in others (for example, M2a and M2b are contained in M2).

4321 Figure B-1 (on page 3420) shows an example of such a model. In a system like this, an  
 4322 application should be able to perform the following operations:

4323



\* All addresses in pool M<sub>2</sub> (comprising pools M<sub>2a</sub> and M<sub>2b</sub>) accessible via port B<sub>1</sub>.  
 Addresses in pool M<sub>2b</sub> are also accessible via port B<sub>2</sub>  
 Addresses in pool M<sub>2a</sub> are NOT accessible via port B<sub>2</sub>

4324

**Figure B-1** Example of a System with Typed Memory

4325

#### — Typed Memory Allocation

4326

4327

4328

4329

4330

4331

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.

4332

#### — Using the Same Storage Region from Different Busses

4333

4334

4335

4336

4337

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.

4338

#### — Sharing Typed Memory Regions

4339

4340

4341

4342

4343

4344

4345

4346

4347

4348

4349

4350

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.

4351

#### — Contiguous Allocation

4352

4353

4354

4355

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



4356 this block of memory. Then, it can pass this offset (and the length) to a process in P2 using  
 4357 some interprocess communication mechanism. The second process can map the same  
 4358 block of memory by using the offset transferred and specifying M3-B2.

4359 — Unallocated Mapping

4360 Any subarea of a memory pool that is mapped to a process, either as the result of an  
 4361 allocation request or an explicit mapping, is normally unavailable for allocation. Special  
 4362 processes such as debuggers, however, may need to map large areas of a typed memory  
 4363 pool, yet leave those areas available for allocation.

4364 Typed memory allocation and mapping has to coexist with storage allocation operators like  
 4365 *malloc()*, but systems are free to choose how to implement this coexistence. For example, it  
 4366 may be system configuration-dependent if all available system memory is made part of one  
 4367 of the typed memory pools or if some part will be restricted to conventional allocation  
 4368 operators. Equally system configuration-dependent may be the availability of operators like  
 4369 *malloc()* to allocate storage from certain typed memory pools. It is not excluded to configure  
 4370 a system such that a given named pool, P1, is in turn split into non-overlapping named  
 4371 subpools. For example, M1-B1, M2-B1, and M3-B1 could also be accessed as one common  
 4372 pool M123-B1. A call to *malloc()* on P1 could work on such a larger pool while full  
 4373 optimization of memory usage by P1 would require typed memory allocation at the subpool  
 4374 level.

4375 • Existing Practice

4376 OS-9 provides for the naming (numbering) and prioritization of memory types by a system  
 4377 administrator. It then provides APIs to request memory allocation of typed (colored)  
 4378 memory by number, and to generate a bus address from a mapped memory address  
 4379 (translate). When requesting colored memory, the user can specify type 0 to signify allocation  
 4380 from the first available type in priority order.

4381 HP-RT presents interfaces to map different kinds of storage regions that are visible through a  
 4382 VME bus, although it does not provide allocation operations. It also provides functions to  
 4383 perform address translation between VME addresses and virtual addresses. It represents a  
 4384 VME-bus unique solution to the general problem.

4385 The PSOS approach is similar (that is, based on a pre-established mapping of bus address  
 4386 ranges to specific memories) with a concept of segments and regions (regions dynamically  
 4387 allocated from a heap which is a special segment). Therefore, PSOS does not fully address the  
 4388 general allocation problem either. PSOS does not have a “process”-based model, but more of  
 4389 a “thread”-only-based model of multi-tasking. So mapping to a process address space is not  
 4390 an issue.

4391 QNX (a Canadian OS vendor specializing in realtime embedded systems on 80x86-based  
 4392 processors) uses the System V approach of opening specially named devices (shared memory  
 4393 segments) and using *mmap()* to then gain access from the process. They do not address  
 4394 allocation directly, but once typed shared memory can be mapped, an “allocation manager”  
 4395 process could be written to handle requests for allocation.

4396 The System V approach also included allocation, implemented by opening yet other special  
 4397 “devices” which allocate, rather than appearing as a whole memory object.

4398 The Orkid realtime kernel interface definition has operations to manage memory “regions”  
 4399 and “pools”, which are areas of memory that may reflect the differing physical nature of the  
 4400 memory. Operations to allocate memory from these regions and pools are also provided.

4401 • Requirements

4402 Existing practice in SVID-derived UNIX systems relies on functionality similar to *mmap()*  
4403 and its related interfaces to achieve mapping and allocation of typed memory. However, the  
4404 issue of sharing typed memory (allocated or mapped) and the complication of multiple ports  
4405 are not addressed in any consistent way by existing UNIX system practice. Part of this  
4406 functionality is existing practice in specialized realtime operating systems. In order to  
4407 solidify the capabilities implied by the model above, the following requirements are imposed  
4408 on the interface:

4409 — Identification of Typed Memory Pools and Ports

4410 All processes (running in all processors) in the system shall be able to identify a particular  
4411 (system configured) typed memory pool accessed through a particular (system  
4412 configured) port by a name. That name shall be a member of a name space common to all  
4413 these processes, but need not be the same name space as that containing ordinary file  
4414 names. The association between memory pools/ports and corresponding names is  
4415 typically established when the system is configured. The “open” operation for typed  
4416 memory objects should be distinct from the *open()* function, for consistency with other  
4417 similar services, but implementable on top of *open()*. This implies that the handle for a  
4418 typed memory object will be a file descriptor.

4419 — Allocation and Mapping of Typed Memory

4420 Once a typed memory object has been identified by a process, it shall be possible to both  
4421 map user-selected subareas of that object into process address space and to map system-  
4422 selected (that is, dynamically allocated) subareas of that object, with user-specified  
4423 length, into process address space. It shall also be possible to determine the maximum  
4424 length of memory allocation that may be requested from a given typed memory object.

4425 — Sharing Typed Memory

4426 Two or more processes shall be able to share portions of typed memory, either user-  
4427 selected or dynamically allocated. This requirement applies also to dynamically allocated  
4428 regions of memory that are composed of several non-contiguous pieces.

4429 — Contiguous Allocation

4430 For dynamic allocation, it shall be the user’s option whether the system is required to  
4431 allocate a contiguous subarea within the typed memory object, or whether it is permitted  
4432 to allocate discontinuous fragments which appear contiguous in the process mapping.  
4433 Contiguous allocation simplifies the process of sharing allocated typed memory, while  
4434 discontinuous allocation allows for potentially better recovery of deallocated typed  
4435 memory.

4436 — Accessing Typed Memory Through Different Ports

4437 Once a subarea of a typed memory object has been mapped, it shall be possible to  
4438 determine the location and length corresponding to a user-selected portion of that object  
4439 within the memory pool. This location and length can then be used to remap that portion  
4440 of memory for access from another port. If the referenced portion of typed memory was  
4441 allocated discontinuously, the length thus determined may be shorter than anticipated,  
4442 and the user code shall adapt to the value returned.

4443 — Deallocation

4444 When a previously mapped subarea of typed memory is no longer mapped by any  
4445 process in the system—as a result of a call or calls to *munmap()*—that subarea shall  
4446 become potentially reusable for dynamic allocation; actual reuse of the subarea is a  
4447 function of the dynamic typed memory allocation policy.

- 4448 — Unallocated Mapping
- 4449 It shall be possible to map user-selected subareas of a typed memory object without  
4450 marking that subarea as unavailable for allocation. This option is not the default behavior,  
4451 and shall require appropriate privilege.
- 4452 • Scenario
- 4453 The following scenario will serve to clarify the use of the typed memory interfaces.
- 4454 Process A running on P1 (see Figure B-1 (on page 3420)) wants to allocate some memory  
4455 from memory pool M2, and it wants to share this portion of memory with process B running  
4456 on P2. Since P2 only has access to the lower part of M2, both processes will use the memory  
4457 pool named M2b which is the part of M2 that is accessible both from P1 and P2. The  
4458 operations that both processes need to perform are shown below:
- 4459 — Allocating Typed Memory
- 4460 Process A calls *posix\_typed\_mem\_open()* with the name **/typed.m2b-b1** and a *tflag* of  
4461 POSIX\_TYPED\_MEM\_ALLOCATE to get a file descriptor usable for allocating from pool  
4462 M2b accessed through port B1. It then calls *mmap()* with this file descriptor requesting a  
4463 length of 4 096 bytes. The system allocates two discontinuous blocks of sizes 1 024 and  
4464 3 072 bytes within M2b. The *mmap()* function returns a pointer to a 4 096 byte array in  
4465 process A's logical address space, mapping the allocated blocks contiguously. Process A  
4466 can then utilize the array, and store data in it.
- 4467 — Determining the Location of the Allocated Blocks
- 4468 Process A can determine the lengths and offsets (relative to M2b) of the two blocks  
4469 allocated, by using the following procedure: First, process A calls *posix\_mem\_offset()*  
4470 with the address of the first element of the array and length 4 096. Upon return, the offset and  
4471 length (1 024 bytes) of the first block are returned. A second call to *posix\_mem\_offset()* is  
4472 then made using the address of the first element of the array plus 1 024 (the length of the  
4473 first block), and a new length of 4 096–1 024. If there were more fragments allocated, this  
4474 procedure could have been continued within a loop until the offsets and lengths of all the  
4475 blocks were obtained. Notice that this relatively complex procedure can be avoided if  
4476 contiguous allocation is requested (by opening the typed memory object with the *tflag*  
4477 POSIX\_TYPED\_MEM\_ALLOCATE\_CONTIG).
- 4478 — Sharing Data Across Processes
- 4479 Process A passes the two offset values and lengths obtained from the *posix\_mem\_offset()*  
4480 calls to process B running on P2, via some form of interprocess communication. Process B  
4481 can gain access to process A's data by calling *posix\_typed\_mem\_open()* with the name  
4482 **/typed.m2b-b2** and a *tflag* of zero, then using two *mmap()* calls on the resulting file  
4483 descriptor to map the two subareas of that typed memory object to its own address space.
- 4484 • Rationale for no *mem\_alloc()* and *mem\_free()*
- 4485 The standard developers had originally proposed a pair of new flags to *mmap()* which, when  
4486 applied to a typed memory object descriptor, would cause *mmap()* to allocate dynamically  
4487 from an unallocated and unmapped area of the typed memory object. Deallocation was  
4488 similarly accomplished through the use of *munmap()*. This was rejected by the ballot group  
4489 because it excessively complicated the (already rather complex) *mmap()* interface and  
4490 introduced semantics useful only for typed memory, to a function which must also map  
4491 shared memory and files. They felt that a memory allocator should be built on top of *mmap()*  
4492 instead of being incorporated within the same interface, much as the ISO C standard libraries  
4493 build *malloc()* on top of the virtual memory mapping functions *brk()* and *sbrk()*. This would

4494 eliminate the complicated semantics involved with unmapping only part of an allocated  
4495 block of typed memory.

4496 To attempt to achieve ballot group consensus, typed memory allocation and deallocation was  
4497 first migrated from *mmap()* and *munmap()* to a pair of complementary functions modeled on  
4498 the ISO C standard *malloc()* and *free()*. The *mem\_alloc()* function specified explicitly the  
4499 typed memory object (typed memory pool/access port) from which allocation takes place,  
4500 unlike *malloc()* where the memory pool and port are unspecified. The *mem\_free()* function  
4501 handled deallocation. These new semantics still met all of the requirements detailed above  
4502 without modifying the behavior of *mmap()* except to allow it to map specified areas of typed  
4503 memory objects. An implementation would have been free to implement *mem\_alloc()* and  
4504 *mem\_free()* over *mmap()*, through *mmap()*, or independently but cooperating with *mmap()*.

4505 The ballot group was queried to see if this was an acceptable alternative, and while there was  
4506 some agreement that it achieved the goal of removing the complicated semantics of  
4507 allocation from the *mmap()* interface, several balloters realized that it just created two  
4508 additional functions that behaved, in great part, like *mmap()*. These balloters proposed an  
4509 alternative which has been implemented here in place of a separate *mem\_alloc()* and  
4510 *mem\_free()*. This alternative is based on four specific suggestions:

- 4511 1. The *posix\_typed\_mem\_open()* function should provide a flag which specifies “allocate  
4512 on *mmap()*” (otherwise, *mmap()* just maps the underlying object). This allows things  
4513 roughly similar to */dev/zero* versus */dev/swap*. Two such flags have been implemented,  
4514 one of which forces contiguous allocation.
- 4515 2. The *posix\_mem\_offset()* function is acceptable because it can be applied usefully to  
4516 mapped objects in general. It should return the file descriptor of the underlying object.
- 4517 3. The *mem\_get\_info()* function in an earlier draft should be renamed  
4518 *posix\_typed\_mem\_get\_info()* because it is not generally applicable to memory objects. It  
4519 should probably return the file descriptor’s allocation attribute. We have implemented  
4520 the renaming of the function, but reject having it return a piece of information which is  
4521 readily known by an application without this function. Its whole purpose is to query  
4522 the typed memory object for attributes that are not user-specified, but determined by  
4523 the implementation.
- 4524 4. There should be no separate *mem\_alloc()* or *mem\_free()* functions. Instead, using  
4525 *mmap()* on a typed memory object opened with an “allocate on *mmap()*” flag should be  
4526 used to force allocation. These are precisely the semantics defined in the current draft.

4527 • Rationale for no Typed Memory Access Management

4528 The working group had originally defined an additional interface (and an additional kind of  
4529 object: typed memory master) to establish and dissolve mappings to typed memory on  
4530 behalf of devices or processors which were independent of the operating system and had no  
4531 inherent capability to directly establish mappings on their own. This was to have provided  
4532 functionality similar to device driver interfaces such as *physio()* and their underlying bus-  
4533 specific interfaces (for example, *mballoc()*) which serve to set up and break down DMA  
4534 pathways, and derive mapped addresses for use by hardware devices and processor cards.

4535 The ballot group felt that this was beyond the scope of POSIX.1 and its amendments.  
4536 Furthermore, the removal of interrupt handling interfaces from a preceding amendment (the  
4537 IEEE Std. 1003.1d-1999) during its balloting process renders these typed memory access  
4538 management interfaces an incomplete solution to portable device management from a user  
4539 process; it would be possible to initiate a device transfer to/from typed memory, but  
4540 impossible to handle the transfer-complete interrupt in a portable way.

4541 To achieve ballot group consensus, all references to typed memory access management  
4542 capabilities were removed. The concept of portable interfaces from a device driver to both  
4543 operating system and hardware is being addressed by the Uniform Driver Interface (UDI)  
4544 industry forum, with formal standardization deferred until proof of concept and industry-  
4545 wide acceptance and implementation.

#### 4546 *B.2.8.4 Process Scheduling*

4547 This portion of the rationale presents models, requirements, and standardization issues relevant  
4548 to process scheduling; see also Section B.2.9.4 (on page 3464).

4549 In an operating system supporting multiple concurrent processes, the system determines the  
4550 order in which processes execute to meet system-defined goals. For time-sharing systems, the  
4551 goal is to enhance system throughput and promote fairness; the application is provided little or  
4552 no control over this sequencing function. While this is acceptable and desirable behavior in a  
4553 time-sharing system, it is inappropriate in a realtime system; realtime applications must  
4554 specifically control the execution sequence of their concurrent processes in order to meet  
4555 externally defined response requirements.

4556 In IEEE Std. 1003.1-200x, the control over process sequencing is provided using a concept of  
4557 scheduling policies. These policies, described in detail in this section, define the behavior of the  
4558 system whenever processor resources are to be allocated to competing processes. Only the  
4559 behavior of the policy is defined; conforming implementations are free to use any mechanism  
4560 desired to achieve the described behavior.

#### 4561 • Models

4562 In an operating system supporting multiple concurrent processes, the system determines the  
4563 order in which processes execute and might force long-running processes to yield to other  
4564 processes at certain intervals. Typically, the scheduling code is executed whenever an event  
4565 occurs that might alter the process to be executed next.

4566 The simplest scheduling strategy is a “first-in, first-out” (FIFO) dispatcher. Whenever a  
4567 process becomes runnable, it is placed on the end of a ready list. The process at the front of  
4568 the ready list is executed until it exits or becomes blocked, at which point it is removed from  
4569 the list. This scheduling technique is also known as “run-to-completion” or “run-to-block”.

4570 A natural extension to this scheduling technique is the assignment of a “non-migrating  
4571 priority” to each process. This policy differs from strict FIFO scheduling in only one respect:  
4572 whenever a process becomes runnable, it is placed at the end of the list of processes runnable  
4573 at that priority level. When selecting a process to run, the system always selects the first  
4574 process from the highest priority queue with a runnable process. Thus, when a process  
4575 becomes unblocked, it will preempt a running process of lower priority without otherwise  
4576 altering the ready list. Further, if a process elects to alter its priority, it is removed from the  
4577 ready list and reinserted, using its new priority, according to the policy above.

4578 While the above policy might be considered unfriendly in a time-sharing environment in  
4579 which multiple users require more balanced resource allocation, it could be ideal in a  
4580 realtime environment for several reasons. The most important of these is that it is  
4581 deterministic: the highest-priority process is always run and, among processes of equal  
4582 priority, the process that has been runnable for the longest time is executed first. Because of  
4583 this determinism, cooperating processes can implement more complex scheduling simply by  
4584 altering their priority. For instance, if processes at a single priority were to reschedule  
4585 themselves at fixed time intervals, a time-slice policy would result.

4586 In a dedicated operating system in which all processes are well-behaved realtime  
4587 applications, non-migrating priority scheduling is sufficient. However, many existing

4588 implementations provide for more complex scheduling policies.

4589 IEEE Std. 1003.1-200x specifies a linear scheduling model. In this model, every process in the  
4590 system has a priority. The system scheduler always dispatches a process that has the highest  
4591 (generally the most time-critical) priority among all runnable processes in the system. As  
4592 long as there is only one such process, the dispatching policy is trivial. When multiple  
4593 processes of equal priority are eligible to run, they are ordered according to a strict run-to-  
4594 completion (FIFO) policy.

4595 The priority is represented as a positive integer and is inherited from the parent process. For  
4596 processes running under a fixed priority scheduling policy, the priority is never altered  
4597 except by an explicit function call.

4598 It was determined arbitrarily that larger integers correspond to “higher priorities”.

4599 Certain implementations might impose restrictions on the priority ranges to which processes  
4600 can be assigned. There also can be restrictions on the set of policies to which processes can be  
4601 set.

4602 • Requirements

4603 Realtime processes require that scheduling be fast and deterministic, and that it guarantees  
4604 to preempt lower priority processes.

4605 Thus, given the linear scheduling model, realtime processes require that they be run at a  
4606 priority that is higher than other processes. Within this framework, realtime processes are  
4607 free to yield execution resources to each other in a completely portable and implementation-  
4608 defined manner.

4609 As there is a generally perceived requirement for processes at the same priority level to share  
4610 processor resources more equitably, provisions are made by providing a scheduling policy  
4611 (that is, SCHED\_RR) intended to provide a timeslice-like facility.

4612 **Note:** The following topics assume that low numeric priority implies low scheduling  
4613 criticality and *vice versa*.

4614 • Rationale for New Interface

4615 Realtime applications need to be able to determine when processes will run in relation to  
4616 each other. It must be possible to guarantee that a critical process will run whenever it is  
4617 runnable; that is, whenever it wants to for as long as it needs. SCHED\_FIFO satisfies this  
4618 requirement. Additionally, SCHED\_RR was defined to meet a realtime requirement for a  
4619 well-defined time-sharing policy for processes at the same priority.

4620 It would be possible to use the BSD *setpriority()* and *getpriority()* functions by redefining the  
4621 meaning of the “nice” parameter according to the scheduling policy currently in use by the  
4622 process. The System V *nice()* interface was felt to be undesirable for realtime because it  
4623 specifies an adjustment to the “nice” value, rather than setting it to an explicit value.  
4624 Realtime applications will usually want to set priority to an explicit value. Also, System V  
4625 *nice()* does not allow for changing the priority of another process.

4626 With the POSIX.1b interfaces, the traditional “nice” value does not affect the SCHED\_FIFO  
4627 or SCHED\_RR scheduling policies. If a “nice” value is supported, it is implementation-  
4628 defined whether it affects the SCHED\_OTHER policy.

4629 An important aspect of IEEE Std. 1003.1-200x is the explicit description of the queuing and  
4630 preemption rules. It is critical, to achieve deterministic scheduling, that such rules be stated  
4631 clearly in IEEE Std. 1003.1-200x.

4632 IEEE Std. 1003.1-200x does not address the interaction between priority and swapping. The  
4633 issues involved with swapping and virtual memory paging are extremely implementation-  
4634 defined and would be nearly impossible to standardize at this point. The proposed  
4635 scheduling paradigm, however, fully describes the scheduling behavior of runnable  
4636 processes, of which one criterion is that the working set be resident in memory. Assuming  
4637 the existence of a portable interface for locking portions of a process in memory, paging  
4638 behavior need not affect the scheduling of realtime processes.

4639 IEEE Std. 1003.1-200x also does not address the priorities of “system” processes. In general,  
4640 these processes should always execute in low-priority ranges to avoid conflict with other  
4641 realtime processes. Implementations should document the priority ranges in which system  
4642 processes run.

4643 The default scheduling policy is not defined. The effect of I/O interrupts and other system  
4644 processing activities is not defined. The temporary lending of priority from one process to  
4645 another (such as for the purposes of affecting freeing resources) by the system is not  
4646 addressed. Preemption of resources is not addressed. Restrictions on the ability of a process  
4647 to affect other processes beyond a certain level (influence levels) is not addressed.

4648 The rationale used to justify the simple time-quantum scheduler is that it is common practice  
4649 to depend upon this type of scheduling to assure “fair” distribution of processor resources  
4650 among portions of the application that must interoperate in a serial fashion. Note that  
4651 IEEE Std. 1003.1-200x is silent with respect to the setting of this time quantum, or whether it  
4652 is a system-wide value or a per-process value, although it appears that the prevailing  
4653 realtime practice is for it to be a system-wide value.

4654 In a system with  $N$  processes at a given priority, all processor-bound, in which the time  
4655 quantum is equal for all processes at a specific priority level, the following assumptions are  
4656 made of such a scheduling policy:

- 4657 1. A time quantum  $Q$  exists and the current process will own control of the processor for  
4658 at least a duration of  $Q$  and will have the processor for a duration of  $Q$ .
- 4659 2. The  $N$ th process at that priority will control a processor within a duration of  $(N-1) \times Q$ .

4660 These assumptions are necessary to provide equal access to the processor and bounded  
4661 response from the application.

4662 The assumptions hold for the described scheduling policy only if no system overhead, such  
4663 as interrupt servicing, is present. If the interrupt servicing load is non-zero, then one of the  
4664 two assumptions becomes fallacious, based upon how  $Q$  is measured by the system.

4665 If  $Q$  is measured by clock time, then the assumption that the process obtains a duration  $Q$   
4666 processor time is false if interrupt overhead exists. Indeed, a scenario can be constructed with  
4667  $N$  processes in which a single process undergoes complete processor starvation if a  
4668 peripheral device, such as an analog-to-digital converter, generates significant interrupt  
4669 activity periodically with a period of  $N \times Q$ .

4670 If  $Q$  is measured as actual processor time, then the assumption that the  $N$ th process runs in  
4671 within the duration  $(N-1) \times Q$  is false.

4672 It should be noted that SCHED\_FIFO suffers from interrupt-based delay as well. However,  
4673 for SCHED\_FIFO, the implied response of the system is “as soon as possible”, so that the  
4674 interrupt load for this case is a vendor selection and not a compliance issue.

4675 With this in mind, it is necessary either to complete the definition by including bounds on the  
4676 interrupt load, or to modify the assumptions that can be made about the scheduling policy.

4677 Since the motivation of inclusion of the policy is common usage, and since current  
4678 applications do not enjoy the luxury of bounded interrupt load, item (2) above is sufficient to  
4679 express existing application needs and is less restrictive in the standard definition. No  
4680 difference in interface is necessary.

4681 In an implementation in which the time quantum is equal for all processes at a specific  
4682 priority, our assumptions can then be restated as:

4683 — A time quantum  $Q$  exists, and a processor-bound process will be rescheduled after a  
4684 duration of, at most,  $Q$ . Time quantum  $Q$  may be defined in either wall clock time or  
4685 execution time.

4686 — In general, the  $N$ th process of a priority level should wait no longer than  $(N-1) \times Q$  time  
4687 to execute, assuming no processes exist at higher priority levels.

4688 — No process should wait indefinitely.

4689 For implementations supporting per-process time quanta, these assumptions can be readily  
4690 extended.

### 4691 **Sporadic Server Scheduling Policy**

4692 The sporadic server is a mechanism defined for scheduling aperiodic activities in time-critical  
4693 realtime systems. This mechanism reserves a certain bounded amount of execution capacity for  
4694 processing aperiodic events at a high priority level. Any aperiodic events that cannot be  
4695 processed within the bounded amount of execution capacity are executed in the background at a  
4696 low priority level. Thus, a certain amount of execution capacity can be guaranteed to be  
4697 available for processing periodic tasks, even under burst conditions in the arrival of aperiodic  
4698 processing requests (that is, a large number of requests in a short time interval). The sporadic  
4699 server also simplifies the schedulability analysis of the realtime system, because it allows  
4700 aperiodic processes or threads to be treated as if they were periodic. The sporadic server was  
4701 first described by Sprunt, et al.

4702 The key concept of the sporadic server is to provide and limit a certain amount of computation  
4703 capacity for processing aperiodic events at their assigned normal priority, during a time interval  
4704 called the *replenishment period*. Once the entity controlled by the sporadic server mechanism is  
4705 initialized with its period and execution-time budget attributes, it preserves its execution  
4706 capacity until an aperiodic request arrives. The request will be serviced (if there are no higher  
4707 priority activities pending) as long as there is execution capacity left. If the request is completed,  
4708 the actual execution time used to service it is subtracted from the capacity, and a replenishment  
4709 of this amount of execution time is scheduled to happen one replenishment period after the  
4710 arrival of the aperiodic request. If the request is not completed, because there is no execution  
4711 capacity left, then the aperiodic process or thread is assigned a lower background priority. For  
4712 each portion of consumed execution capacity the execution time used is replenished after one  
4713 replenishment period. At the time of replenishment, if the sporadic server was executing at a  
4714 background priority level, its priority is elevated to the normal level. Other similar  
4715 replenishment policies have been defined, but the one presented here represents a compromise  
4716 between efficiency and implementation complexity.

4717 The interface that appears in this section defines a new scheduling policy for threads and  
4718 processes that behaves according to the rules of the sporadic server mechanism. Scheduling  
4719 attributes are defined and functions are provided to allow the user to set and get the parameters  
4720 that control the scheduling behavior of this mechanism, namely the normal and low priority, the  
4721 replenishment period, the maximum number of pending replenishment operations, and the  
4722 initial execution-time budget.



- 4723 • Scheduling Aperiodic Activities
- 4724 Virtually all realtime applications are required to process aperiodic activities. In many cases,  
4725 there are tight timing constraints that the response to the aperiodic events must meet. Usual  
4726 timing requirements imposed on the response to these events are:
  - 4727 — The effects of an aperiodic activity on the response time of lower priority activities must  
4728 be controllable and predictable.
  - 4729 — The system must provide the fastest possible response time to aperiodic events.
  - 4730 — It must be possible to take advantage of all the available processing bandwidth not  
4731 needed by time-critical activities to enhance average-case response times to aperiodic  
4732 events.
- 4733 Traditional methods for scheduling aperiodic activities are background processing, polling  
4734 tasks, and direct event execution:
  - 4735 — Background processing consists of assigning a very low priority to the processing of  
4736 aperiodic events. It utilizes all the available bandwidth in the system that has not been  
4737 consumed by higher priority threads. However, it is very difficult, or impossible, to meet  
4738 requirements on average-case response time, because the aperiodic entity has to wait for  
4739 the execution of all other entities which have higher priority.
  - 4740 — Polling consists of creating a periodic process or thread for servicing aperiodic requests.  
4741 At regular intervals, the polling entity is started and it services accumulated pending  
4742 aperiodic requests. If no aperiodic requests are pending, the polling entity suspends itself  
4743 until its next period. Polling allows the aperiodic requests to be processed at a higher  
4744 priority level. However, worst and average-case response times of polling entities are a  
4745 direct function of the polling period, and there is execution overhead for each polling  
4746 period, even if no event has arrived. If the deadline of the aperiodic activity is short  
4747 compared to the inter-arrival time, the polling frequency must be increased to guarantee  
4748 meeting the deadline. For this case, the increase in frequency can dramatically reduce the  
4749 efficiency of the system and, therefore, its capacity to meet all deadlines. Yet, polling  
4750 represents a good way to handle a large class of practical problems because it preserves  
4751 system predictability, and because the amortized overhead drops as load increases.
  - 4752 — Direct event execution consists of executing the aperiodic events at a high fixed-priority  
4753 level. Typically, the aperiodic event is processed by an interrupt service routine as soon as  
4754 it arrives. This technique provides predictable response times for aperiodic events, but  
4755 makes the response times of all lower priority activities completely unpredictable under  
4756 burst arrival conditions. Therefore, if the density of aperiodic event arrivals is  
4757 unbounded, it may be a dangerous technique for time-critical systems. Yet, for those cases  
4758 in which the physics of the system imposes a bound on the event arrival rate, it is  
4759 probably the most efficient technique.
  - 4760 — The sporadic server scheduling algorithm combines the predictability of the polling  
4761 approach with the short response times of the direct event execution. Thus, it allows  
4762 systems to meet an important class of application requirements that cannot be met by  
4763 using the traditional approaches. Multiple sporadic servers with different attributes can  
4764 be applied to the scheduling of multiple classes of aperiodic events, each with different  
4765 kinds of timing requirements, such as individual deadlines, average response times, and  
4766 so on. It also has many other interesting applications for realtime, such as scheduling  
4767 producer/consumer tasks in time-critical systems, limiting the effects of faults on the  
4768 estimation of task execution-time requirements, and so on.

- 4769       • Existing Practice  
4770       The sporadic server has been used in different kinds of applications, including military  
4771       avionics, robot control systems, industrial automation systems, and so on. There are  
4772       examples of many systems that cannot be successfully scheduled using the classic  
4773       approaches, such as direct event execution, or polling, and are schedulable using a sporadic  
4774       server scheduler. The sporadic server algorithm itself can successfully schedule all systems  
4775       scheduled with direct event execution or polling.  
  
4776       The sporadic server scheduling policy has been implemented as a commercial product in the  
4777       run-time system of the Verdex Ada compiler. There are also many applications that have  
4778       used a much less efficient application-level sporadic server. These real-time applications  
4779       would benefit from a sporadic server scheduler implemented at the scheduler level.
- 4780       • Library-Level *versus* Kernel-Level Implementation  
  
4781       The sporadic server interface described in this section requires the sporadic server policy to  
4782       be implemented at the same level as the scheduler. This means that the process sporadic  
4783       server shall be implemented at the kernel level and the thread sporadic server policy shall be  
4784       implemented at the same level as the thread scheduler; that is, kernel or library level.  
  
4785       In an earlier interface for the sporadic server, this mechanism was implementable at a  
4786       different level than the scheduler. This feature allowed the implementer to choose between  
4787       an efficient scheduler-level implementation, or a simpler user or library-level  
4788       implementation. However, the working group considered that this interface made the use of  
4789       sporadic servers more complex, and that library-level implementations would lack some of  
4790       the important functionality of the sporadic server, namely the limitation of the actual  
4791       execution time of aperiodic activities. The working group also felt that the interface  
4792       described in this chapter does not preclude library-level implementations of threads intended  
4793       to provide efficient low-overhead scheduling for those threads that are not scheduled under  
4794       the sporadic server policy.
- 4795       • Range of Scheduling Priorities  
  
4796       Each of the scheduling policies supported in IEEE Std. 1003.1-200x has an associated range of  
4797       priorities. The priority ranges for each policy might or might not overlap with the priority  
4798       ranges of other policies. For time-critical realtime applications it is usual for periodic and  
4799       aperiodic activities to be scheduled together in the same processor. Periodic activities will  
4800       usually be scheduled using the SCHED\_FIFO scheduling policy, while aperiodic activities  
4801       may be scheduled using SCHED\_SPORADIC. Since the application developer will require  
4802       complete control over the relative priorities of these activities in order to meet his timing  
4803       requirements, it would be desirable for the priority ranges of SCHED\_FIFO and  
4804       SCHED\_SPORADIC to overlap completely. Therefore, although IEEE Std. 1003.1-200x does  
4805       not require any particular relationship between the different priority ranges, it is  
4806       recommended that these two ranges should coincide.
- 4807       • Dynamically Setting the Sporadic Server Policy  
  
4808       Several members of the working group requested that implementations should not be  
4809       required to support dynamically setting the sporadic server scheduling policy for a thread.  
4810       The reason is that this policy may have a high overhead for library-level implementations of  
4811       threads, and if threads are allowed to dynamically set this policy, this overhead can be  
4812       experienced even if the thread does not use that policy. By disallowing the dynamic setting  
4813       of the sporadic server scheduling policy, these implementations can accomplish efficient  
4814       scheduling for threads using other policies. If a strictly conforming application needs to use  
4815       the sporadic server policy, and is therefore willing to pay the overhead, it must set this policy  
4816       at the time of thread creation.

4817 • Limitation of the Number of Pending Replenishments

4818 The number of simultaneously pending replenishment operations must be limited for each  
 4819 sporadic server for two reasons: an unlimited number of replenishment operations would  
 4820 need an unlimited number of system resources to store all the pending replenishment  
 4821 operations; on the other hand, in some implementations each replenishment operation will  
 4822 represent a source of priority inversion (just for the duration of the replenishment operation)  
 4823 and thus, the maximum amount of replenishments must be bounded to guarantee bounded  
 4824 response times. The way in which the number of replenishments is bounded is by lowering  
 4825 the priority of the sporadic server to *sched\_ss\_low\_priority* when the number of pending  
 4826 replenishments has reached its limit. In this way, no new replenishments are scheduled until  
 4827 the number of pending replenishments decreases.

4828 In the sporadic server scheduling policy defined in IEEE Std. 1003.1-200x, the application can  
 4829 specify the maximum number of pending replenishment operations for a single sporadic  
 4830 server, by setting the value of the *sched\_ss\_max\_repl* scheduling parameter. This value must  
 4831 be between one and {SS\_REPL\_MAX}, which is a maximum limit imposed by the  
 4832 implementation. The limit {SS\_REPL\_MAX} must be greater than or equal to  
 4833 {\_POSIX\_SS\_REPL\_MAX}, which is defined to be four in IEEE Std. 1003.1-200x. The  
 4834 minimum limit of four was chosen so that an application can at least guarantee that four  
 4835 different aperiodic events can be processed during each interval of length equal to the  
 4836 replenishment period.

4837 **B.2.8.5** *Clocks and Timers*

4838 • Clocks

4839 IEEE Std. 1003.1-200x and the ISO C standard both define functions for obtaining system  
 4840 time. Implicit behind these functions is a mechanism for measuring passage of time. This  
 4841 specification makes this mechanism explicit and calls it a clock. The *CLOCK\_REALTIME*  
 4842 clock required by IEEE Std. 1003.1-200x is a higher resolution version of the clock that  
 4843 maintains POSIX.1 system time. This is a “system-wide” clock, in that it is visible to all  
 4844 processes and, were it possible for multiple processes to all read the clock at the same time,  
 4845 they would see the same value.

4846 An extensible interface was defined, with the ability for implementations to define additional  
 4847 clocks. This was done because of the observation that many realtime platforms support  
 4848 multiple clocks, and it was desired to fit this model within the standard interface. But  
 4849 implementation-defined clocks need not represent actual hardware devices, nor are they  
 4850 necessarily system-wide.

4851 • Timers

4852 Two timer types are required for a system to support realtime applications:

4853 1. One-shot

4854 A one-shot timer is a timer that is armed with an initial expiration time, either relative  
 4855 to the current time or at an absolute time (based on some timing base, such as time in  
 4856 seconds and nanoseconds since the Epoch). The timer expires once and then is  
 4857 disarmed. With the specified facilities, this is accomplished by setting the *it\_value*  
 4858 member of the *value* argument to the desired expiration time and the *it\_interval* member  
 4859 to zero.

4860 2. Periodic

4861 A periodic timer is a timer that is armed with an initial expiration time, again either  
 4862 relative or absolute, and a repetition interval. When the initial expiration occurs, the

4863 timer is reloaded with the repetition interval and continues counting. With the  
 4864 specified facilities, this is accomplished by setting the *it\_value* member of the *value*  
 4865 argument to the desired initial expiration time and the *it\_interval* member to the desired  
 4866 repetition interval.

4867 For both of these types of timers, the time of the initial timer expiration can be specified in  
 4868 two ways:

- 4869 1. Relative (to the current time)
- 4870 2. Absolute

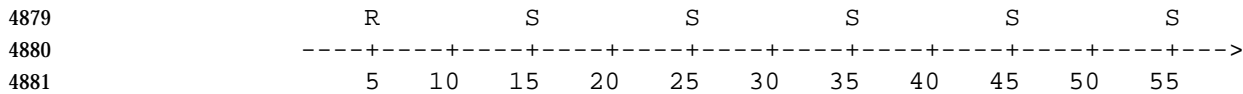
4871 • Examples of Using Realtime Timers

4872 In the diagrams below, *S* indicates a program schedule, *R* shows a schedule method request,  
 4873 and *E* suggests an internal operating system event.

4874 — Periodic Timer: Data Logging

4875 During an experiment, it might be necessary to log realtime data periodically to an  
 4876 internal buffer or to a mass storage device. With a periodic scheduling method, a logging  
 4877 module can be started automatically at fixed time intervals to log the data.

4878 Program schedule is requested every 10 seconds.



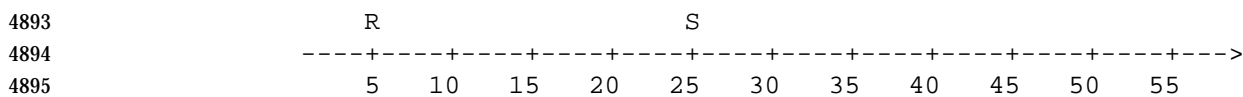
4882 [Time (in Seconds)]

4883 To achieve this type of scheduling using the specified facilities, one would allocate a per-  
 4884 process timer based on clock ID `CLOCK_REALTIME`. Then the timer would be armed via  
 4885 a call to `timer_settime()` with the `TIMER_ABSTIME` flag reset, and with an initial  
 4886 expiration value and a repetition interval of 10 seconds.

4887 — One-shot Timer (Relative Time): Device Initialization

4888 In an emission test environment, large sample bags are used to capture the exhaust from  
 4889 a vehicle. The exhaust is purged from these bags before each and every test. With a one-  
 4890 shot timer, a module could initiate the purge function and then suspend itself for a  
 4891 predetermined period of time while the sample bags are prepared.

4892 Program schedule requested 20 seconds after call is issued.



4896 [Time (in Seconds)]

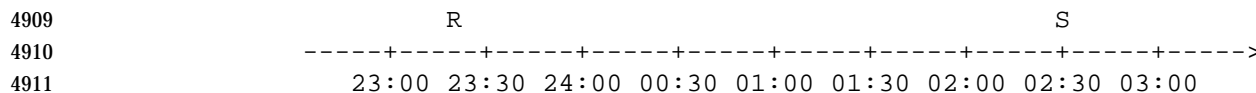
4897 To achieve this type of scheduling using the specified facilities, one would allocate a per-  
 4898 process timer based on clock ID `CLOCK_REALTIME`. Then the timer would be armed via  
 4899 a call to `timer_settime()` with the `TIMER_ABSTIME` flag reset, and with an initial  
 4900 expiration value of 20 seconds and a repetition interval of zero.

4901 Note that if the program wishes merely to suspend itself for the specified interval, it  
 4902 could more easily use `nanosleep()`.

4903 — One-shot Timer (Absolute Time): Data Transmission

4904 The results from an experiment are often moved to a different system within a network  
 4905 for postprocessing or archiving. With an absolute one-shot timer, a module that moves  
 4906 data from a test-cell computer to a host computer can be automatically scheduled on a  
 4907 daily basis.

4908 Program schedule requested for 2:30 a.m.



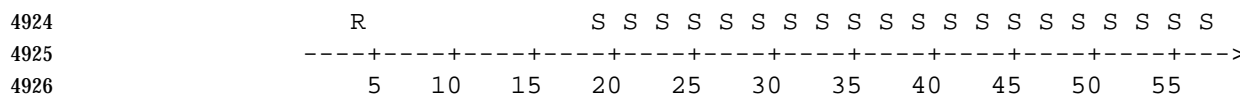
4912 [Time of Day]

4913 To achieve this type of scheduling using the specified facilities, one would allocate a per-  
 4914 process timer based on clock ID CLOCK\_REALTIME. Then the timer would be armed via  
 4915 a call to *timer\_settime()* with the *TIMER\_ABSTIME* flag set, and an initial expiration value  
 4916 equal to 2:30 a.m. of the next day.

4917 — Periodic Timer (Relative Time): Signal Stabilization

4918 Some measurement devices, such as emission analyzers, do not respond instantaneously  
 4919 to an introduced sample. With a periodic timer with a relative initial expiration time, a  
 4920 module that introduces a sample and records the average response could suspend itself  
 4921 for a predetermined period of time while the signal is stabilized and then sample at a  
 4922 fixed rate.

4923 Program schedule requested 15 seconds after call is issued and every 2 seconds thereafter.



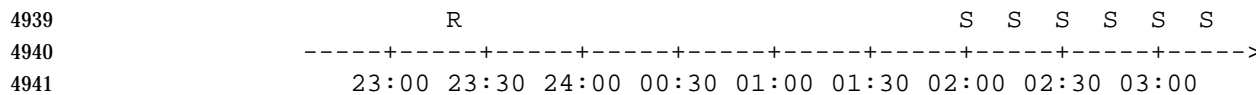
4927 [Time (in Seconds)]

4928 To achieve this type of scheduling using the specified facilities, one would allocate a per-  
 4929 process timer based on clock ID CLOCK\_REALTIME. Then the timer would be armed via  
 4930 a call to *timer\_settime()* with *TIMER\_ABSTIME* flag reset, and with an initial expiration  
 4931 value of 15 seconds and a repetition interval of 2 seconds.

4932 — Periodic Timer (Absolute Time): Work Shift-related Processing

4933 Resource utilization data is useful when time to perform experiments is being scheduled  
 4934 at a facility. With a periodic timer with an absolute initial expiration time, a module can  
 4935 be scheduled at the beginning of a work shift to gather resource utilization data  
 4936 throughout the shift. This data can be used to allocate resources effectively to minimize  
 4937 bottlenecks and delays and maximize facility throughput.

4938 Program schedule requested for 2:00 a.m. and every 15 minutes thereafter.



4942 [Time of Day]

4943 To achieve this type of scheduling using the specified facilities, one would allocate a per-  
 4944 process timer based on clock ID CLOCK\_REALTIME. Then the timer would be armed via  
 4945 a call to *timer\_settime()* with *TIMER\_ABSTIME* flag set, and with an initial expiration  
 4946 value equal to 2:00 a.m. and a repetition interval equal to 15 minutes.

4947 • Relationship of Timers to Clocks

4948 The relationship between clocks and timers armed with an absolute time is straightforward:  
 4949 a timer expiration signal is requested when the associated clock reaches or exceeds the  
 4950 specified time. The relationship between clocks and timers armed with a relative time (an  
 4951 interval) is less obvious, but not unintuitive. In this case, a timer expiration signal is  
 4952 requested when the specified interval, *as measured by the associated clock*, has passed. For the  
 4953 required CLOCK\_REALTIME clock, this allows timer expiration signals to be requested at  
 4954 specified “wall clock” times (absolute), or when a specified interval of “realtime” has passed  
 4955 (relative). For an implementation-defined clock—say, a process virtual time clock—timer  
 4956 expirations could be requested when the process has used a specified total amount of virtual  
 4957 time (absolute), or when it has used a specified *additional* amount of virtual time (relative).

4958 The interfaces also allow flexibility in the implementation of the functions. For example, an  
 4959 implementation could convert all absolute times to intervals by subtracting the clock value at  
 4960 the time of the call from the requested expiration time and “counting down” at the  
 4961 supported resolution. Or it could convert all relative times to absolute expiration time by  
 4962 adding in the clock value at the time of the call and comparing the clock value to the  
 4963 expiration time at the supported resolution. Or it might even choose to maintain absolute  
 4964 times as absolute and compare them to the clock value at the supported resolution for  
 4965 absolute timers, and maintain relative times as intervals and count them down at the  
 4966 resolution supported for relative timers. The choice will be driven by efficiency  
 4967 considerations and the underlying hardware or software clock implementation.

4968 • Data Definitions for Clocks and Timers

4969 IEEE Std. 1003.1-200x uses a time representation capable of supporting nanosecond  
 4970 resolution timers for the following reasons:

- 4971 — To enable IEEE Std. 1003.1-200x to represent those computer systems already using  
 4972 nanosecond or submicrosecond resolution clocks.
- 4973 — To accommodate those per-process timers that might need nanoseconds to specify an  
 4974 absolute value of system-wide clocks, even though the resolution of the per-process timer  
 4975 may only be milliseconds, or *vice versa*.
- 4976 — Because the number of nanoseconds in a second can be represented in 32 bits.

4977 Time values are represented in the **timespec** structure. The *tv\_sec* member is of type **time\_t**  
 4978 so that this member is compatible with time values used by POSIX.1 functions and the ISO C  
 4979 standard. The *tv\_nsec* member is a **signed long** in order to simplify and clarify code that  
 4980 decrements or finds differences of time values. Note that because 1 billion (number of  
 4981 nanoseconds per second) is less than half of the value representable by a signed 32-bit value,  
 4982 it is always possible to add two valid fractional seconds represented as integral nanoseconds  
 4983 without overflowing the signed 32-bit value.

4984 A maximum allowable resolution for the CLOCK\_REALTIME clock of 20 ms (1/50 seconds)  
 4985 was chosen to allow line frequency clocks in European countries to be conforming. 60 Hz  
 4986 clocks in the U.S. will also be conforming, as will finer granularity clocks, although a Strictly  
 4987 Conforming Application cannot assume a granularity of less than 20 ms (1/50 seconds).

4988 The minimum allowable maximum time allowed for the CLOCK\_REALTIME clock and the  
 4989 function *nanosleep()*, and timers created with *clock\_id*=CLOCK\_REALTIME, is determined by  
 4990 the fact that the *tv\_sec* member is of type **time\_t**.

4991 IEEE Std. 1003.1-200x specifies that timer expirations shall not be delivered early, nor shall  
 4992 *nanosleep()* return early due to quantization error. IEEE Std. 1003.1-200x discusses the various  
 4993 implementations of *alarm()* in the rationale and states that implementations that do not

4994 allow alarm signals to occur early are the most appropriate, but refrained from mandating  
4995 this behavior. Because of the importance of predictability to realtime applications,  
4996 IEEE Std. 1003.1-200x takes a stronger stance.

4997 The developers of IEEE Std. 1003.1-200x considered using a time representation that differs  
4998 from POSIX.1b in the second 32 bit of the 64-bit value. Whereas POSIX.1b defines this field  
4999 as a fractional second in nanoseconds, the other methodology defines this as a binary fraction  
5000 of one second, with the radix point assumed before the most significant bit.

5001 POSIX.1b is a software, source-level standard and most of the benefits of the alternate  
5002 representation are enjoyed by hardware implementations of clocks and algorithms. It was  
5003 felt that mandating this format for POSIX.1b clocks and timers would unnecessarily burden  
5004 the application writer with writing, possibly non-portable, multiple precision arithmetic  
5005 packages to perform conversion between binary fractions and integral units such as  
5006 nanoseconds, milliseconds, and so on.

#### 5007 **Rationale for the Monotonic Clock**

5008 For those applications that use time services to achieve realtime behavior, changing the value of  
5009 the clock on which these services rely may cause erroneous timing behavior. For these  
5010 applications, it is necessary to have a monotonic clock which cannot run backwards, and which  
5011 has a maximum clock jump that is required to be documented by the implementation.  
5012 Additionally, it is desirable (but not required by IEEE Std. 1003.1-200x) that the monotonic clock  
5013 increases its value uniformly. This clock should not be affected by changes to the system time;  
5014 for example, to synchronize the clock with an external source or to account for leap seconds.  
5015 Such changes would cause errors in the measurement of time intervals for those time services  
5016 that use the absolute value of the clock.

5017 One could argue that by defining the behavior of time services when the value of a clock is  
5018 changed, deterministic realtime behavior can be achieved. For example, one could specify that  
5019 relative time services should be unaffected by changes in the value of a clock. However, there  
5020 are time services that are based upon an absolute time, but that are essentially intended as  
5021 relative time services. For example, *pthread\_cond\_timedwait()* uses an absolute time to allow it to  
5022 wake up after the required interval despite spurious wakeups. Although sometimes the  
5023 *pthread\_cond\_timedwait()* timeouts are absolute in nature, there are many occasions in which  
5024 they are relative, and their absolute value is determined from the current time plus a relative  
5025 time interval. In this latter case, if the clock changes while the thread is waiting, the wait interval  
5026 will not be the expected length. If a *pthread\_cond\_timedwait()* function were created that would  
5027 take a relative time, it would not solve the problem because to retain the intended “deadline” a  
5028 thread would need to compensate for latency due to the spurious wakeup, and preemption  
5029 between wakeup and the next wait.

5030 The solution is to create a new monotonic clock, whose value does not change except for the  
5031 regular ticking of the clock, and use this clock for implementing the various relative timeouts  
5032 that appear in the different POSIX interfaces, as well as allow *pthread\_cond\_timedwait()* to choose  
5033 this new clock for its timeout. A new *clock\_nanosleep()* function is created to allow an application  
5034 to take advantage of this newly defined clock. Notice that the monotonic clock may be  
5035 implemented using the same hardware clock as the system clock.

5036 Relative timeouts for *sigtimedwait()* and *aio\_suspend()* have been redefined to use the monotonic  
5037 clock, if present. The *alarm()* function has not been redefined, because the same effect but with  
5038 better resolution can be achieved by creating a timer (for which the appropriate clock may be  
5039 chosen).

5040 The *pthread\_cond\_timedwait()* function has been treated in a different way, compared to other  
5041 functions with absolute timeouts, because it is used to wait for an event, and thus it may have a

5042 deadline, while the other timeouts are generally used as an error recovery mechanism, and for  
5043 them the use of the monotonic clock is not so important. Since the desired timeout for the  
5044 *pthread\_cond\_timedwait()* function may either be a relative interval, or an absolute time of day  
5045 deadline, a new initialization attribute has been created for condition variables, to specify the  
5046 clock that shall be used for measuring the timeout in a call to *pthread\_cond\_timedwait()*. In this  
5047 way, if a relative timeout is desired, the monotonic clock will be used; if an absolute deadline is  
5048 required instead, the `CLOCK_REALTIME` or another appropriate clock may be used. This  
5049 capability has not been added to other functions with absolute timeouts because for those  
5050 functions the expected use of the timeout is mostly to prevent errors, and not so often to meet  
5051 precise deadlines. As a consequence, the complexity of adding this capability is not justified by  
5052 its perceived application usage.

5053 The *nanosleep()* function has not been modified with the introduction of the monotonic clock.  
5054 Instead, a new *clock\_nanosleep()* function has been created, in which the desired clock may be  
5055 specified in the function call.

5056 • History of Resolution Issues

5057 Due to the shift from relative to absolute timeouts in IEEE Std. 1003.1d-1999, the  
5058 amendments to the *sem\_timedwait()*, *pthread\_mutex\_timedlock()*, *mq\_timedreceive()*, and  
5059 *mq\_timedsend()* functions of that standard have been removed. Those amendments specified  
5060 that `CLOCK_MONOTONIC` would be used for the (relative) timeouts if the Monotonic  
5061 Clock option was supported.

5062 Having these functions continue to be tied solely to `CLOCK_MONOTONIC` would not  
5063 work. Since the absolute value of a time value obtained from `CLOCK_MONOTONIC` is  
5064 unspecified, under the absolute timeouts interface, applications would behave differently  
5065 depending on whether the Monotonic Clock option was supported or not (because the  
5066 absolute value of the clock would have different meanings in either case).

5067 Two options were considered:

- 5068 1. Leave the current behavior unchanged, which specifies the `CLOCK_REALTIME` clock  
5069 for these (absolute) timeouts, to allow portability of applications between  
5070 implementations supporting or not the Monotonic Clock option.
- 5071 2. Modify these functions in the way that *pthread\_cond\_timedwait()* was modified to allow  
5072 a choice of clock, so that an application could use `CLOCK_REALTIME` when it is trying  
5073 to achieve an absolute timeout and `CLOCK_MONOTONIC` when it is trying to achieve  
5074 a relative timeout.

5075 It was decided that the features of `CLOCK_MONOTONIC` are not as critical to these  
5076 functions as they are to *pthread\_cond\_timedwait()*. The *pthread\_cond\_timedwait()* function is  
5077 given a relative timeout; the timeout may represent a deadline for an event. When these  
5078 functions are given relative timeouts, the timeouts are typically for error recovery purposes  
5079 and need not be so precise.

5080 Therefore, it was decided that these functions should be tied to `CLOCK_REALTIME` and not  
5081 complicated by being given a choice of clock.



5082           **Execution Time Monitoring**

## 5083           • Introduction

5084           The main goals of the execution time monitoring facilities defined in this chapter are to  
5085           measure the execution time of processes and threads and to allow an application to establish  
5086           CPU time limits for these entities.

5087           The analysis phase of time-critical realtime systems often relies on the measurement of  
5088           execution times of individual threads or processes to determine whether the timing  
5089           requirements will be met. Also, performance analysis techniques for soft deadline realtime  
5090           systems rely heavily on the determination of these execution times. The execution time  
5091           monitoring functions provide application developers with the ability to measure these  
5092           execution times online and open the possibility of dynamic execution-time analysis and  
5093           system reconfiguration, if required.

5094           The second goal of allowing an application to establish execution time limits for individual  
5095           processes or threads and detecting when they overrun allows program robustness to be  
5096           increased by enabling online checking of the execution times.

5097           If errors are detected—possibly because of erroneous program constructs, the existence of  
5098           errors in the analysis phase, or a burst of event arrivals—online detection and recovery is  
5099           possible in a portable way. This feature can be extremely important for many time-critical  
5100           applications. Other applications require trapping CPU-time errors as a normal way to exit an  
5101           algorithm; for instance, some realtime artificial intelligence applications trigger a number of  
5102           independent inference processes of varying accuracy and speed, limit how long they can run,  
5103           and pick the best answer available when time runs out. In many periodic systems, overrun  
5104           processes are simply restarted in the next resource period, after necessary end-of-period  
5105           actions have been taken. This allows algorithms that are inherently data-dependent to be  
5106           made predictable.

5107           The interface that appears in this chapter defines a new type of clock, the CPU-time clock,  
5108           which measures execution time. Each process or thread can invoke the clock and timer  
5109           functions defined in POSIX.1 to use them. Functions are also provided to access the CPU-  
5110           time clock of other processes or threads to enable remote monitoring of these clocks.  
5111           Monitoring of threads of other processes is not supported, since these threads are not visible  
5112           from outside of their own process with the interfaces defined in POSIX.1.

## 5113           • Execution Time Monitoring Interface

5114           The clock and timer interface defined in POSIX.1 historically only defined one clock, which  
5115           measures wall-clock time. The requirements for measuring execution time of processes and  
5116           threads, and setting limits to their execution time by detecting when they overrun, can be  
5117           accomplished with that interface if a new kind of clock is defined. These new clocks measure  
5118           execution time, and one is associated with each process and with each thread. The clock  
5119           functions currently defined in POSIX.1 can be used to read and set these CPU-time clocks,  
5120           and timers can be created using these clocks as their timing base. These timers can then be  
5121           used to send a signal when some specified execution time has been exceeded. The CPU-time  
5122           clocks of each process or thread can be accessed by using the symbols  
5123           CLOCK\_PROCESS\_CPUTIME\_ID or CLOCK\_THREAD\_CPUTIME\_ID.

5124           The clock and timer interface defined in POSIX.1 and extended with the new kind of CPU-  
5125           time clock would only allow processes or threads to access their own CPU-time clocks.  
5126           However, many realtime systems require the possibility of monitoring the execution time of  
5127           processes or threads from independent monitoring entities. In order to allow applications to  
5128           construct independent monitoring entities that do not require cooperation from or  
5129           modification of the monitored entities, two functions have been added: *clock\_getcpu* and *clockid*(),

5130 for accessing CPU-time clocks of other processes, and *pthread\_getcpuclockid()*, for accessing  
5131 CPU-time clocks of other threads. These functions return the clock identifier associated with  
5132 the process or thread specified in the call. These clock IDs can then be used in the rest of the  
5133 clock function calls.

5134 The clocks accessed through these functions could also be used as a timing base for the  
5135 creation of timers, thereby allowing independent monitoring entities to limit the CPU-time  
5136 consumed by other entities. However, this possibility would imply additional complexity  
5137 and overhead because of the need to maintain a timer queue for each process or thread, to  
5138 store the different expiration times associated with timers created by different processes or  
5139 threads. The working group decided this additional overhead was not justified by  
5140 application requirements. Therefore, creation of timers attached to the CPU-time clocks of  
5141 other processes or threads has been specified as implementation-defined.

5142 • Overhead Considerations

5143 The measurement of execution time may introduce additional overhead in the thread  
5144 scheduling, because of the need to keep track of the time consumed by each of these entities.  
5145 In library-level implementations of threads, the efficiency of scheduling could be somehow  
5146 compromised because of the need to make a kernel call, at each context switch, to read the  
5147 process CPU-time clock. Consequently, a thread creation attribute called *cpu-clock-*  
5148 *requirement* was defined, to allow threads to disconnect their respective CPU-time clocks.  
5149 However, the Ballot Group considered that this attribute itself introduced some overhead,  
5150 and that in current implementations it was not worth the effort. Therefore, the attribute was  
5151 deleted, and thus thread CPU-time clocks are required for all threads if the Thread CPU-Time  
5152 Clocks option is supported.

5153 • Accuracy of CPU-time Clocks

5154 The mechanism used to measure the execution time of processes and threads is specified in  
5155 IEEE Std. 1003.1-200x as implementation-defined. The reason for this is that both the  
5156 underlying hardware and the implementation architecture have a very strong influence on  
5157 the accuracy achievable for measuring CPU time. For some implementations, the  
5158 specification of strict accuracy requirements would represent very large overheads, or even  
5159 the impossibility of being implemented.

5160 Since the mechanism for measuring execution time is implementation-defined, realtime  
5161 applications will be able to take advantage of accurate implementations using a portable  
5162 interface. Of course, strictly conforming applications cannot rely on any particular degree of  
5163 accuracy, in the same way as they cannot rely on a very accurate measurement of wall clock  
5164 time. There will always exist applications whose accuracy or efficiency requirements on the  
5165 implementation are more rigid than the values defined in IEEE Std. 1003.1-200x or any other  
5166 standard.

5167 In any case, there is a minimum set of characteristics that realtime applications would expect  
5168 from most implementations. One such characteristic is that the sum of all the execution times  
5169 of all the threads in a process equals the process execution time, when no CPU-time clocks  
5170 are disabled. This need not always be the case because implementations may differ in how  
5171 they account for time during context switches. Another characteristic is that the sum of the  
5172 execution times of all processes in a system equals the number of processors, multiplied by  
5173 the elapsed time, assuming that no processor is idle during that elapsed time. However, in  
5174 some systems it might not be possible to relate CPU-time to elapsed time. For example, in a  
5175 heterogeneous multi-processor system in which each processor runs at a different speed, an  
5176 implementation may choose to define each “second” of CPU-time to be a certain number of  
5177 “cycles” that a CPU has executed.

- 5178           • Existing Practice
- 5179           Measuring and limiting the execution time of each concurrent activity are common features  
5180           of most industrial implementations of realtime systems. Almost all critical realtime systems  
5181           are currently built upon a cyclic executive. With this approach, a regular timer interrupt kicks  
5182           off the next sequence of computations. It also checks that the current sequence has  
5183           completed. If it has not, then some error recovery action can be undertaken (or at least an  
5184           overrun is avoided). Current software engineering principles and the increasing complexity  
5185           of software are driving application developers to implement these systems on multi-  
5186           threaded or multi-process operating systems. Therefore, if a POSIX operating system is to be  
5187           used for this type of application, then it must offer the same level of protection.
- 5188           Execution time clocks are also common in most UNIX implementations, although these  
5189           clocks usually have requirements different from those of realtime applications. The POSIX.1  
5190           *times()* function supports the measurement of the execution time of the calling process, and  
5191           its terminated child processes. This execution time is measured in clock ticks and is supplied  
5192           as two different values with the user and system execution times, respectively. BSD supports  
5193           the function *getrusage()*, which allows the calling process to get information about the  
5194           resources used by itself and/or all of its terminated child processes. The resource usage  
5195           includes user and system CPU time. Some UNIX systems have options to specify high  
5196           resolution (up to one microsecond) CPU time clocks using the *times()* or the *getrusage()*  
5197           functions.
- 5198           The *times()* and *getrusage()* interfaces do not meet important realtime requirements, such as  
5199           the possibility of monitoring execution time from a different process or thread, or the  
5200           possibility of detecting an execution time overrun. The latter requirement is supported in  
5201           some UNIX implementations that are able to send a signal when the execution time of a  
5202           process has exceeded some specified value. For example, BSD defines the functions  
5203           *getitimer()* and *setitimer()*, which can operate either on a realtime clock (wall-clock), or on  
5204           virtual-time or profile-time clocks which measure CPU time in two different ways. These  
5205           functions do not support access to the execution time of other processes.
- 5206           IBM's MVS operating system supports per-process and per-thread execution time clocks. It  
5207           also supports limiting the execution time of a given process.
- 5208           Given all this existing practice, the working group considered that the POSIX.1 clocks and  
5209           timers interface was appropriate to meet most of the requirements that realtime applications  
5210           have for execution time clocks. Functions were added to get the CPU time clock IDs, and to  
5211           allow/disallow the thread CPU time clocks (in order to preserve the efficiency of some  
5212           implementations of threads).
- 5213           • Clock Constants
- 5214           The definition of the manifest constants `CLOCK_PROCESS_CPUTIME_ID` and  
5215           `CLOCK_THREAD_CPUTIME_ID` allows processes or threads, respectively, to access their  
5216           own execution-time clocks. However, given a process or thread, access to its own execution-  
5217           time clock is also possible if the clock ID of this clock is obtained through a call to  
5218           *clock\_getcpuclockid()* or *pthread\_getcpuclockid()*. Therefore, these constants are not necessary  
5219           and could be deleted to make the interface simpler. Their existence saves one system call in  
5220           the first access to the CPU-time clock of each process or thread. The working group  
5221           considered this issue and decided to leave the constants in IEEE Std. 1003.1-200x because  
5222           they are closer to the POSIX.1b use of clock identifiers.
- 5223           • Library Implementations of Threads
- 5224           In library implementations of threads, kernel entities and library threads can coexist. In this  
5225           case, if the CPU-time clocks are supported, most of the clock and timer functions will need to

5226 have two implementations: one in the thread library, and one in the system calls library. The  
 5227 main difference between these two implementations is that the thread library  
 5228 implementation will have to deal with clocks and timers that reside in the thread space,  
 5229 while the kernel implementation will operate on timers and clocks that reside in kernel space.  
 5230 In the library implementation, if the clock ID refers to a clock that resides in the kernel, a  
 5231 kernel call will have to be made. The correct version of the function can be chosen by  
 5232 specifying the appropriate order for the libraries during the link process.

5233 • History of Resolution Issues: Deletion of the *enable* Attribute

5234 In the draft corresponding to the first balloting round, CPU-time clocks had an attribute  
 5235 called *enable*. This attribute was introduced by the working group to allow implementations  
 5236 to avoid the overhead of measuring execution time for those processes or threads for which  
 5237 this measurement was not required. However, the *enable* attribute got several ballot  
 5238 objections. The main reason was that processes are already required to measure execution  
 5239 time by the POSIX.1 *times()* function. Consequently, the *enable* attribute was considered  
 5240 unnecessary, and was deleted from the draft.

#### 5241 Rationale Relating to Timeouts

5242 • Requirements for Timeouts

5243 Realtime systems which must operate reliably over extended periods without human  
 5244 intervention are characteristic in embedded applications such as avionics, machine control,  
 5245 and space exploration, as well as more mundane applications such as cable TV, security  
 5246 systems, and plant automation. A multi-tasking paradigm, in which many independent  
 5247 and/or cooperating software functions relinquish the processor(s) while waiting for a  
 5248 specific stimulus, resource, condition, or operation completion, is very useful in producing  
 5249 well engineered programs for such systems. For such systems to be robust and fault-tolerant,  
 5250 expected occurrences that are unduly delayed or that never occur must be detected so that  
 5251 appropriate recovery actions may be taken. This is difficult if there is no way for a task to  
 5252 regain control of a processor once it has relinquished control (blocked) awaiting an  
 5253 occurrence which, perhaps because of corrupted code, hardware malfunction, or latent  
 5254 software bugs, will not happen when expected. Therefore, the common practice in realtime  
 5255 operating systems is to provide a capability to timeout such blocking services. Although  
 5256 there are several methods to achieve this already defined by POSIX, none are as reliable or  
 5257 efficient as initiating a timeout simultaneously with initiating a blocking service. This is  
 5258 especially critical in hard-realtime embedded systems because the processors typically have  
 5259 little time reserve, and allowed fault recovery times are measured in milliseconds rather than  
 5260 seconds.

5261 The working group largely agreed that such timeouts were necessary and ought to become  
 5262 part of IEEE Std. 1003.1-200x, particularly vendors of realtime operating systems whose  
 5263 customers had already expressed a strong need for timeouts. There was some resistance to  
 5264 inclusion of timeouts in IEEE Std. 1003.1-200x because the desired effect, fault tolerance,  
 5265 could, in theory, be achieved using existing facilities and alternative software designs, but  
 5266 there was no compelling evidence that realtime system designers would embrace such  
 5267 designs at the sacrifice of performance and/or simplicity.

5268 • Which Services should be Timed Out?

5269 Originally, the working group considered the prospect of providing timeouts on all blocking  
 5270 services, including those currently existing in POSIX.1, POSIX.1b, and POSIX.1c, and future  
 5271 interfaces to be defined by other working groups, as sort of a general policy. This was rather  
 5272 quickly rejected because of the scope of such a change, and the fact that many of those  
 5273 services would not normally be used in a realtime context. More traditional timesharing

5274 solutions to timeout would suffice for most of the POSIX.1 interfaces, while others had  
 5275 asynchronous alternatives which, while more complex to utilize, would be adequate for  
 5276 some realtime and all non-realtime applications.

5277 The list of potential candidates for timeouts was narrowed to the following for further  
 5278 consideration:

5279 — POSIX.1b

5280 — *sem\_wait()*

5281 — *mq\_receive()*

5282 — *mq\_send()*

5283 — *lio\_listio()*

5284 — *aio\_suspend()*

5285 — *sigwait()* (timeout already implemented by *sigtimedwait()*)

5286 — POSIX.1c

5287 — *pthread\_mutex\_lock()*

5288 — *pthread\_join()*

5289 — *pthread\_cond\_wait()* (timeout already implemented by *pthread\_cond\_timedwait()*)

5290 — POSIX.1

5291 — *read()*

5292 — *write()*

5293 After further review by the working group, the *lio\_listio()*, *read()*, and *write()* functions (all  
 5294 forms of blocking synchronous I/O) were eliminated from the list because of the following:

5295 — Asynchronous alternatives exist

5296 — Timeouts can be implemented, albeit non-portably, in device drivers

5297 — A strong desire not to introduce modifications to POSIX.1 interfaces

5298 The working group ultimately rejected *pthread\_join()* since both that interface and a timed  
 5299 variant of that interface are non-minimal and may be implemented as a library function. See  
 5300 below for a library implementation of *pthread\_join()*.

5301 Thus, there was a consensus among the working group members to add timeouts to 4 of the  
 5302 remaining 5 functions (the timeout for *aio\_suspend()* was ultimately added directly to  
 5303 POSIX.1b, while the others were added by POSIX.1d). However, *pthread\_mutex\_lock()*  
 5304 remained contentious.

5305 Many feel that *pthread\_mutex\_lock()* falls into the same class as the other functions; that is, it  
 5306 is desirable to timeout a mutex lock because a mutex may fail to be unlocked due to errant or  
 5307 corrupted code in a critical section (looping or branching outside of the unlock code), and  
 5308 therefore is equally in need of a reliable, simple, and efficient timeout. In fact, since mutexes  
 5309 are intended to guard small critical sections, most *pthread\_mutex\_lock()* calls would be  
 5310 expected to obtain the lock without blocking nor utilizing any kernel service, even in  
 5311 implementations of threads with global contention scope; the timeout alternative need only  
 5312 be considered after it is determined that the thread must block.

5313 Those opposed to timing out mutexes feel that the very simplicity of the mutex is  
 5314 compromised by adding a timeout semantic, and that to do so is senseless. They claim that if

5315 a timed mutex is really deemed useful by a particular application, then it can be constructed  
 5316 from the facilities already in POSIX.1b and POSIX.1c. The following two C-language library  
 5317 implementations of mutex locking with timeout represent the solutions offered (in both  
 5318 implementations, the timeout parameter is specified as absolute time, not relative time as in  
 5319 the proposed POSIX.1c interfaces).

5320 • Spinlock Implementation

```

5321     #include <pthread.h>
5322     #include <time.h>
5323     #include <errno.h>

5324     int pthread_mutex_timedlock(pthread_mutex_t *mutex,
5325                               const struct timespec *timeout)
5326     {
5327         struct timespec timenow;

5328         while (pthread_mutex_trylock(mutex) == EBUSY)
5329             {
5330                 clock_gettime(CLOCK_REALTIME, &timenow);
5331                 if (timespec_cmp(&timenow, timeout) >= 0)
5332                     {
5333                         return ETIMEDOUT;
5334                     }
5335                 pthread_yield();
5336             }
5337         return 0;
5338     }
  
```

5339 The Spinlock implementation is generally unsuitable for any application using priority-based  
 5340 thread scheduling policies such as SCHED\_FIFO or SCHED\_RR, since the mutex could  
 5341 currently be held by a thread of lower priority within the same allocation domain, but since  
 5342 the waiting thread never blocks, only threads of equal or higher priority will ever run, and  
 5343 the mutex cannot be unlocked. Setting priority inheritance or priority ceiling protocol on the  
 5344 mutex does not solve this problem, since the priority of a mutex owning thread is only  
 5345 boosted if higher priority threads are blocked waiting for the mutex; clearly not the case for  
 5346 this spinlock.

5347 • Condition Wait Implementation

```

5348     #include <pthread.h>
5349     #include <time.h>
5350     #include <errno.h>

5351     struct timed_mutex
5352     {
5353         int locked;
5354         pthread_mutex_t mutex;
5355         pthread_cond_t cond;
5356     };
5357     typedef struct timed_mutex timed_mutex_t;

5358     int timed_mutex_lock(timed_mutex_t *tm,
5359                       const struct timespec *timeout)
5360     {
5361         int timedout=FALSE;
5362         int error_status;
  
```

```

5363     pthread_mutex_lock(&tm->mutex);
5364     while (tm->locked && !timedout)
5365     {
5366         if ((error_status=pthread_cond_timedwait(&tm->cond,
5367         &tm->mutex,
5368         timeout))!=0)
5369         {
5370             if (error_status==ETIMEDOUT) timedout = TRUE;
5371         }
5372     }
5373     if(timedout)
5374     {
5375         pthread_mutex_unlock(&tm->mutex);
5376         return ETIMEDOUT;
5377     }
5378     else
5379     {
5380         tm->locked = TRUE;
5381         pthread_mutex_unlock(&tm->mutex);
5382         return 0;
5383     }
5384 }
5385 void timed_mutex_unlock(timed_mutex_t *tm)
5386 {
5387     pthread_mutex_lock(&tm->mutex); / for case assignment not atomic /
5388     tm->locked = FALSE;
5389     pthread_mutex_unlock(&tm->mutex);
5390     pthread_cond_signal(&tm->cond);
5391 }

```

5392 The Condition Wait implementation effectively substitutes the *pthread\_cond\_timedwait()*  
5393 function (which is currently timed out) for the desired *pthread\_mutex\_timedlock()*. Since waits  
5394 on condition variables currently do not include protocols which avoid priority inversion, this  
5395 method is generally unsuitable for realtime applications because it does not provide the same  
5396 priority inversion protection as the untimed *pthread\_mutex\_lock()*. Also, for any given  
5397 implementations of the current mutex and condition variable primitives, this library  
5398 implementation has a performance cost at least 2.5 times that of the untimed  
5399 *pthread\_mutex\_lock()* even in the case where the timed mutex is readily locked without  
5400 blocking (the interfaces required for this case are shown in bold). Even in uniprocessors or  
5401 where assignment is atomic, at least an additional *pthread\_cond\_signal()* is required.  
5402 *pthread\_mutex\_timedlock()* could be implemented at effectively no performance penalty in  
5403 this case because the timeout parameters need only be considered after it is determined that  
5404 the mutex cannot be locked immediately.

5405 Thus it has not yet been shown that the full semantics of mutex locking with timeout can be  
5406 efficiently and reliably achieved using existing interfaces. Even if the existence of an  
5407 acceptable library implementation were proven, it is difficult to justify why the interface  
5408 itself should not be made portable, especially considering approval for the other four  
5409 timeouts.

5410 • Rationale for Library Implementation of *pthread\_timedjoin()*

```

5411     Library implementation of pthread_timedjoin():
5412     /*
5413     * Construct a thread variety entirely from existing functions
5414     * with which a join can be done, allowing the join to time out.
5415     */
5416     #include <pthread.h>
5417     #include <time.h>
5418
5418     struct timed_thread {
5419         pthread_t t;
5420         pthread_mutex_t m;
5421         int exiting;
5422         pthread_cond_t exit_c;
5423         void *(*start_routine)(void *arg);
5424         void *arg;
5425         void *status;
5426     };
5427
5427     typedef struct timed_thread *timed_thread_t;
5428     static pthread_key_t timed_thread_key;
5429     static pthread_once_t timed_thread_once = PTHREAD_ONCE_INIT;
5430
5430     static void timed_thread_init()
5431     {
5432         pthread_key_create(&timed_thread_key, NULL);
5433     }
5434
5434     static void *timed_thread_start_routine(void *args)
5435     /*
5436     * Routine to establish thread-specific data value and run the actual
5437     * thread start routine which was supplied to timed_thread_create().
5438     */
5439     {
5440         timed_thread_t tt = (timed_thread_t) args;
5441
5441         pthread_once(&timed_thread_once, timed_thread_init);
5442         pthread_setspecific(timed_thread_key, (void *)tt);
5443         timed_thread_exit((tt->start_routine)(tt->arg));
5444     }
5445
5445     int timed_thread_create(timed_thread_t ttp, const pthread_attr_t *attr,
5446         void *(*start_routine)(void *), void *arg)
5447     /*
5448     * Allocate a thread which can be used with timed_thread_join().
5449     */
5450     {
5451         timed_thread_t tt;
5452         int result;
5453
5453         tt = (timed_thread_t) malloc(sizeof(struct timed_thread));
5454         pthread_mutex_init(&tt->m, NULL);
5455         tt->exiting = FALSE;
5456         pthread_cond_init(&tt->exit_c, NULL);
5457         tt->start_routine = start_routine;

```



```

5458         tt->arg = arg;
5459         tt->status = NULL;

5460         if ((result = pthread_create(&tt->t, attr,
5461             timed_thread_start_routine, (void *)tt)) != 0) {
5462             free(tt);
5463             return result;
5464         }

5465         pthread_detach(tt->t);
5466         ttp = tt;
5467         return 0;
5468     }

5469     int timed_thread_join(timed_thread_t tt,
5470         struct timespec *timeout,
5471         void **status)
5472     {
5473         int result;

5474         pthread_mutex_lock(&tt->m);
5475         result = 0;
5476         /*
5477          * Wait until the thread announces that it is exiting,
5478          * or until timeout.
5479          */
5480         while (result == 0 && ! tt->exiting) {
5481             result = pthread_cond_timedwait(&tt->exit_c, &tt->m, timeout);
5482         }
5483         pthread_mutex_unlock(&tt->m);
5484         if (result == 0 && tt->exiting) {
5485             *status = tt->status;
5486             free((void *)tt);
5487             return result;
5488         }
5489         return result;
5490     }

5491     void timed_thread_exit(void *status)
5492     {
5493         timed_thread_t tt;
5494         void *specific;

5495         if ((specific=pthread_getspecific(timed_thread_key)) == NULL){
5496             /*
5497              * Handle cases which won't happen with correct usage.
5498              */
5499             pthread_exit( NULL);
5500         }
5501         tt = (timed_thread_t) specific;
5502         pthread_mutex_lock(&tt->m);
5503         /*
5504          * Tell a joiner that we're exiting.
5505          */
5506         tt->status = status;

```

```

5507         tt->exiting = TRUE;
5508         pthread_cond_signal(&tt->exit_c);
5509         pthread_mutex_unlock(&tt->m);
5510         /*
5511          * Call pthread_exit() to call destructors and really
5512          * exit the thread.
5513          */
5514         pthread_exit(NULL);
5515     }

```

5516 The *pthread\_join()* C-language example shown above demonstrates that it is possible, using  
5517 existing pthread facilities, to construct a variety of thread which allows for joining such a  
5518 thread, but which allows the join operation to time out. It does this by using a  
5519 *pthread\_cond\_timedwait()* to wait for the thread to exit. A **timed\_thread\_t** descriptor structure  
5520 is used to pass parameters from the creating thread to the created thread, and from the  
5521 exiting thread to the joining thread. This implementation is roughly equivalent to what a  
5522 normal *pthread\_join()* implementation would do, with the single change being that  
5523 *pthread\_cond\_timedwait()* is used in place of a simple *pthread\_cond\_wait()*.

5524 Since it is possible to implement such a facility entirely from existing pthread interfaces, and  
5525 with roughly equal efficiency and complexity to an implementation which would be  
5526 provided directly by a pthreads implementation, it was the consensus of the working group  
5527 members that any *pthread\_timedjoin()* facility would be unnecessary, and should not be  
5528 provided.

#### 5529 • Form of the Timeout Interfaces

5530 The working group considered a number of alternative ways to add timeouts to blocking  
5531 services. At first, a system interface which would specify a one-shot or persistent timeout to  
5532 be applied to subsequent blocking services invoked by the calling process or thread was  
5533 considered because it allowed all blocking services to be timed out in a uniform manner with  
5534 a single additional interface; this was rather quickly rejected because it could easily result in  
5535 the wrong services being timed out.

5536 It was suggested that a timeout value might be specified as an attribute of the object  
5537 (semaphore, mutex, message queue, and so on), but there was no consensus on this, either on  
5538 a case-by-case basis or for all timeouts.

5539 Looking at the two existing timeouts for blocking services indicates that the working group  
5540 members favor a separate interface for the timed version of a function. However,  
5541 *pthread\_cond\_timedwait()* utilizes an absolute timeout value while *sigtimedwait()* uses a  
5542 relative timeout value. The working group members agreed that relative timeout values are  
5543 appropriate where the timeout mechanism's primary use was to deal with an unexpected or  
5544 error situation, but they are inappropriate when the timeout must expire at a particular time,  
5545 or before a specific deadline. For the timeouts being introduced in IEEE Std. 1003.1-200x, the  
5546 working group considered allowing both relative and absolute timeouts as is done with  
5547 POSIX.1b timers, but ultimately favored the simpler absolute timeout form.

5548 An absolute time measure can be easily implemented on top of an interface that specifies  
5549 relative time, by reading the clock, calculating the difference between the current time and  
5550 the desired wake-up time, and issuing a relative timeout call. But there is a race condition  
5551 with this approach because the thread could be preempted after reading the clock, but before  
5552 making the timed out call; in this case, the thread would be awakened later than it should  
5553 and, thus, if the wake up time represented a deadline, it would miss it.

5554 There is also a race condition when trying to build a relative timeout on top of an interface  
5555 that specifies absolute timeouts. In this case, we would have to read the clock to calculate the  
5556 absolute wake-up time as the sum of the current time plus the relative timeout interval. In  
5557 this case, if the thread is preempted after reading the clock but before making the timed out  
5558 call, the thread would be awakened earlier than desired.

5559 But the race condition with the absolute timeouts interface is not as bad as the one that  
5560 happens with the relative timeout interface, because there are simple workarounds. For the  
5561 absolute timeouts interface, if the timing requirement is a deadline, we can still meet this  
5562 deadline because the thread woke up earlier than the deadline. If the timeout is just used as  
5563 an error recovery mechanism, the precision of timing is not really important. If the timing  
5564 requirement is that between actions A and B a minimum interval of time must elapse, we can  
5565 safely use the absolute timeout interface by reading the clock after action A has been started.  
5566 It could be argued that, since the call with the absolute timeout is atomic from the  
5567 application point of view, it is not possible to read the clock after action A, if this action is  
5568 part of the timed out call. But if we look at the nature of the calls for which we specify  
5569 timeouts (locking a mutex, waiting for a semaphore, waiting for a message, or waiting until  
5570 there is space in a message queue), the timeouts that an application would build on these  
5571 actions would not be triggered by these actions themselves, but by some other external  
5572 action. For example, if we want to wait for a message to arrive to a message queue, and wait  
5573 for at least 20 milliseconds, this time interval would start to be counted from some event that  
5574 would trigger both the action that produces the message, as well as the action that waits for  
5575 the message to arrive, and not by the wait-for-message operation itself. In this case, we could  
5576 use the workaround proposed above.

5577 For these reasons, the absolute timeout is preferred over the relative timeout interface.

## 5578 **B.2.9 Threads**

5579 Threads will normally be more expensive than subroutines (or functions, routines, and so on) if  
5580 specialized hardware support is not provided. Nevertheless, threads should be sufficiently  
5581 efficient to encourage their use as a medium to fine-grained structuring mechanism for  
5582 parallelism in an application. Structuring an application using threads then allows it to take  
5583 immediate advantage of any underlying parallelism available in the host environment. This  
5584 means implementors are encouraged to optimize for fast execution at the possible expense of  
5585 efficient utilization of storage. For example, a common thread creation technique is to cache  
5586 appropriate thread data structures. That is, rather than releasing system resources, the  
5587 implementation retains these resources and reuses them when the program next asks to create a  
5588 new thread. If this reuse of thread resources is to be possible, there has to be very little unique  
5589 state associated with each thread, because any such state has to be reset when the thread is  
5590 reused.

### 5591 **Thread Creation Attributes**

5592 Attributes objects are provided for threads, mutexes, and condition variables as a mechanism to  
5593 support probable future standardization in these areas without requiring that the interface itself  
5594 be changed. Attributes objects provide clean isolation of the configurable aspects of threads. For  
5595 example, “stack size” is an important attribute of a thread, but it cannot be expressed portably.  
5596 When porting a threaded program, stack sizes often need to be adjusted. The use of attributes  
5597 objects can help by allowing the changes to be isolated in a single place, rather than being spread  
5598 across every instance of thread creation.

5599 Attributes objects can be used to set up *classes* of threads with similar attributes; for example,  
5600 “threads with large stacks and high priority” or “threads with minimal stacks”. These classes  
5601 can be defined in a single place and then referenced wherever threads need to be created.

5602 Changes to “class” decisions become straightforward, and detailed analysis of each  
5603 *pthread\_create()* call is not required.

5604 The attributes objects are defined as opaque types as an aid to extensibility. If these objects had  
5605 been specified as structures, adding new attributes would force recompilation of all multi-  
5606 threaded programs when the attributes objects are extended; this might not be possible if  
5607 different program components were supplied by different vendors.

5608 Additionally, opaque attributes objects present opportunities for improving performance.  
5609 Argument validity can be checked once when attributes are set, rather than each time a thread is  
5610 created. Implementations will often need to cache kernel objects that are expensive to create.  
5611 Opaque attributes objects provide an efficient mechanism to detect when cached objects become  
5612 invalid due to attribute changes.

5613 Because assignment is not necessarily defined on a given opaque type, implementation-  
5614 dependent default values cannot be defined in a portable way. The solution to this problem is to  
5615 allow attribute objects to be initialized dynamically by attributes object initialization functions,  
5616 so that default values can be supplied automatically by the implementation.

5617 The following proposal was provided as a suggested alternative to the supplied attributes:

- 5618 1. Maintain the style of passing a parameter formed by the bitwise-inclusive OR of flags to  
5619 the initialization routines (*pthread\_create()*, *pthread\_mutex\_init()*, *pthread\_cond\_init()*). The  
5620 parameter containing the flags should be an opaque type for extensibility. If no flags are  
5621 set in the parameter, then the objects are created with default characteristics. An  
5622 implementation may specify implementation-defined flag values and associated behavior.
- 5623 2. If further specialization of mutexes and condition variables is necessary, implementations  
5624 may specify additional procedures that operate on the **pthread\_mutex\_t** and  
5625 **pthread\_cond\_t** objects (instead of on attributes objects).

5626 The difficulties with this solution are:

- 5627 1. A bitmask is not opaque if bits have to be set into bit-vector attributes objects using  
5628 explicitly-coded bitwise-inclusive OR operations. If the set of options exceeds an **int**,  
5629 application programmers need to know the location of each bit. If bits are set or read by  
5630 encapsulation (that is, *get\*()* or *set\*()* functions), then the bitmask is merely an  
5631 implementation of attributes objects as currently defined and should not be exposed to the  
5632 programmer.
- 5633 2. Many attributes are not Boolean or very small integral values. For example, scheduling  
5634 policy may be placed in 3 bits or 4 bits, but priority requires 5 bits or more, thereby taking  
5635 up at least 8 bits out of a possible 16 bits on machines with 16-bit integers. Because of this,  
5636 the bitmask can only reasonably control whether particular attributes are set or not, and it  
5637 cannot serve as the repository of the value itself. The value needs to be specified as a  
5638 function parameter (which is non-extensible), or by setting a structure field (which is non-  
5639 opaque), or by *get\*()* and *set\*()* functions (making the bitmask a redundant addition to the  
5640 attributes objects).

5641 Stack size is defined as an optional attribute because the very notion of a stack is inherently  
5642 machine-dependent. Some implementations may not be able to change the size of the stack, for  
5643 example, and others may not need to because stack pages may be discontinuous and can be  
5644 allocated and released on demand.

5645 The attribute mechanism has been designed in large measure for extensibility. Future extensions  
5646 to the attribute mechanism or to any attributes object defined in IEEE Std. 1003.1-200x has to be  
5647 done with care so as not to affect binary-compatibility.

5648 Attribute objects, even if allocated by means of dynamic allocation functions such as *malloc()*,  
5649 may have their size fixed at compile time. This means, for example, a *pthread\_create()* in an  
5650 implementation with extensions to the **pthread\_attr\_t** cannot look beyond the area that the  
5651 binary application assumes is valid. This suggests that implementations should maintain a size  
5652 field in the attributes object, as well as possibly version information, if extensions in different  
5653 directions (possibly by different vendors) are to be accommodated.

#### 5654 **Thread Implementation Models**

5655 There are various thread implementation models. At one end of the spectrum is the “library-  
5656 thread model”. In such a model, the threads of a process are not visible to the operating system  
5657 kernel, and the threads are not kernel scheduled entities. The process is the only kernel  
5658 scheduled entity. The process is scheduled onto the processor by the kernel according to the  
5659 scheduling attributes of the process. The threads are scheduled onto the single kernel scheduled  
5660 entity (the process) by the runtime library according to the scheduling attributes of the threads.  
5661 A problem with this model is that it constrains concurrency. Since there is only one kernel  
5662 scheduled entity (namely, the process), only one thread per process can execute at a time. If the  
5663 thread that is executing blocks on I/O, then the whole process blocks.

5664 At the other end of the spectrum is the “kernel-thread model”. In this model, all threads are  
5665 visible to the operating system kernel. Thus, all threads are kernel scheduled entities, and all  
5666 threads can concurrently execute. The threads are scheduled onto processors by the kernel  
5667 according to the scheduling attributes of the threads. The drawback to this model is that the  
5668 creation and management of the threads entails operating system calls, as opposed to subroutine  
5669 calls, which makes kernel threads heavier weight than library threads.

5670 Hybrids of these two models are common. A hybrid model offers the speed of library threads  
5671 and the concurrency of kernel threads. In hybrid models, a process has some (relatively small)  
5672 number of kernel scheduled entities associated with it. It also has a potentially much larger  
5673 number of library threads associated with it. Some library threads may be bound to kernel  
5674 scheduled entities, while the other library threads are multiplexed onto the remaining kernel  
5675 scheduled entities. There are two levels of thread scheduling:

- 5676 1. The runtime library manages the scheduling of (unbound) library threads onto kernel  
5677 scheduled entities.
- 5678 2. The kernel manages the scheduling of kernel scheduled entities onto processors.

5679 For this reason, a hybrid model is referred to as a *two-level threads scheduling model*. In this model,  
5680 the process can have multiple concurrently executing threads; specifically, it can have as many  
5681 concurrently executing threads as it has kernel scheduled entities.

#### 5682 **Thread-Specific Data**

5683 Many applications require that a certain amount of context be maintained on a per-thread basis  
5684 across procedure calls. A common example is a multi-threaded library routine that allocates  
5685 resources from a common pool and maintains an active resource list for each thread. The  
5686 thread-specific data interface provided to meet these needs may be viewed as a two-dimensional  
5687 array of values with keys serving as the row index and thread IDs as the column index (although  
5688 the implementation need not work this way).

##### 5689 • Models

5690 Three possible thread-specific data models were considered:

- 5691 1. No Explicit Support

5692 A standard thread-specific data interface is not strictly necessary to support  
5693 applications that require per-thread context. One could, for example, provide a hash  
5694 function that converted a **pthread\_t** into an integer value that could then be used to  
5695 index into a global array of per-thread data pointers. This hash function, in conjunction  
5696 with *pthread\_self()*, would be all the interface required to support a mechanism of this  
5697 sort. Unfortunately, this technique is cumbersome. It can lead to duplicated code as  
5698 each set of cooperating modules implements their own per-thread data management  
5699 schemes.

## 5700 2. Single (**void \***) Pointer

5701 Another technique would be to provide a single word of per-thread storage and a pair  
5702 of functions to fetch and store the value of this word. The word could then hold a  
5703 pointer to a block of per-thread memory. The allocation, partitioning, and general use  
5704 of this memory would be entirely up to the application. Although this method is not as  
5705 problematic as technique 1, it suffers from interoperability problems. For example, all  
5706 modules using the per-thread pointer would have to agree on a common usage  
5707 protocol.

## 5708 3. Key/Value Mechanism

5709 This method associates an opaque key (for example, stored in a variable of type  
5710 **pthread\_key\_t**) with each per-thread datum. These keys play the role of identifiers for  
5711 per-thread data. This technique is the most generic and avoids the problems noted  
5712 above, albeit at the cost of some complexity.

5713 The primary advantage of the third model is its information hiding properties. Modules  
5714 using this model are free to create and use their own key(s) independent of all other such  
5715 usage, whereas the other models require that all modules that use thread-specific context  
5716 explicitly cooperate with all other such modules. The data-independence provided by the  
5717 third model is worth the additional interface.

### 5718 • Requirements

5719 It is important that it be possible to implement the thread-specific data interface without the  
5720 use of thread private memory. To do otherwise would increase the weight of each thread,  
5721 thereby limiting the range of applications for which the threads interfaces provided by  
5722 IEEE Std. 1003.1-200x is appropriate.

5723 The values that one binds to the key via *pthread\_setspecific()* may, in fact, be pointers to  
5724 shared storage locations available to all threads. It is only the key/value bindings that are  
5725 maintained on a per-thread basis, and these can be kept in any portion of the address space  
5726 that is reserved for use by the calling thread (for example, on the stack). Thus, no per-thread  
5727 MMU state is required to implement the interface. On the other hand, there is nothing in the  
5728 interface specification to preclude the use of a per-thread MMU state if it is available (for  
5729 example, the key values returned by *pthread\_key\_create()* could be thread private memory  
5730 addresses).

### 5731 • Standardization Issues

5732 Thread-specific data is a requirement for a usable thread interface. The binding described in  
5733 this section provides a portable thread-specific data mechanism for languages that do not  
5734 directly support a thread-specific storage class. A binding to IEEE Std. 1003.1-200x for a  
5735 language that does include such a storage class need not provide this specific interface.

5736 If a language were to include the notion of thread-specific storage, it would be desirable (but  
5737 *not* required) to provide an implementation of the pthreads thread-specific data interface  
5738 based on the language feature. For example, assume that a compiler for a C-like language

5739 supports a *private* storage class that provides thread-specific storage. Something similar to  
 5740 the following macros might be used to effect a compatible implementation:

```
5741     #define pthread_key_t                private void *
5742     #define pthread_key_create(key)      /* no-op */
5743     #define pthread_setspecific(key,value) (key)=(value)
5744     #define pthread_getspecific(key)     (key)
```

5745 **Note:** For the sake of clarity, this example ignores destructor functions. A correct  
 5746 implementation would have to support them.

## 5747 **Barriers**

### 5748 • Background

5749 Barriers are typically used in parallel DO/FOR loops to ensure that all threads have reached  
 5750 a particular stage in a parallel computation before allowing any to proceed to the next stage.  
 5751 Highly efficient implementation is possible on machines which support a “Fetch and Add”  
 5752 operation as described in the referenced Almasi and Gottlieb (1989).

5753 The use of return value PTHREAD\_BARRIER\_SERIAL\_THREAD is shown in the following  
 5754 example:

```
5755     if ( (status=pthread_barrier_wait(&barrier)) ==
5756         PTHREAD_BARRIER_SERIAL_THREAD) {
5757         ...serial section
5758     }
5759         else if (status != 0) {
5760         ...error processing
5761     }
5762     status=pthread_barrier_wait(&barrier);
5763     ...
```

5764 This behavior allows a serial section of code to be executed by one thread as soon as all  
 5765 threads reach the first barrier. The second barrier prevents the other threads from proceeding  
 5766 until the serial section being executed by the one thread has completed.

5767 Although barriers can be implemented with mutexes and condition variables, the referenced  
 5768 Almasi and Gottlieb (1989) provides ample illustration that such implementations are  
 5769 significantly less efficient than is possible. While the relative efficiency of barriers may well  
 5770 vary by implementation, it is important that they be recognized in the IEEE Std. 1003.1-200x  
 5771 to facilitate application portability while providing the necessary freedom to implementors.

### 5772 • Lack of Timeout Feature

5773 Alternate versions of most blocking routines have been provided to support watchdog  
 5774 timeouts. No alternate interface of this sort has been provided for barrier waits for the  
 5775 following reasons:

- 5776 • Multiple threads may use different timeout values, some of which may be indefinite. It is  
 5777 not clear which threads should break through the barrier with a timeout error if and when  
 5778 these timeouts expire.
- 5779 • The barrier may become unusable once a thread breaks out of a *pthread\_barrier\_wait()*  
 5780 with a timeout error. There is, in general, no way to guarantee the consistency of a  
 5781 barrier's internal data structures once a thread has timed out of a *pthread\_barrier\_wait()*.  
 5782 Even the inclusion of a special barrier reinitialization function would not help much since  
 5783 it is not clear how this function would affect the behavior of threads that reach the barrier

5784                    between the original timeout and the call to the reinitialization function.

5785                    **Spin Locks**

5786                    • Background

5787                    Spin locks represent an extremely low-level synchronization mechanism suitable primarily  
5788                    for use on shared memory multi-processors. It is typically an atomically modified Boolean  
5789                    value that is set to one when the lock is held and to zero when the lock is freed.

5790                    When a caller requests a spin lock that is already held, it typically spins in a loop testing  
5791                    whether the lock has become available. Such spinning wastes processor cycles so the lock  
5792                    should only be held for short durations and not across sleep/block operations. Callers should  
5793                    unlock spin locks before calling sleep operations.

5794                    Spin locks are available on a variety of systems. The functions included in  
5795                    IEEE Std. 1003.1-200x are an attempt to standardize that existing practice.

5796                    • Lack of Timeout Feature

5797                    Alternate versions of most blocking routines have been provided to support watchdog  
5798                    timeouts. No alternate interface of this sort has been provided for spin locks for the following  
5799                    reasons:

5800                    • It is impossible to determine appropriate timeout intervals for spin locks in a portable  
5801                    manner. The amount of time one can expect to spend spin-waiting is inversely  
5802                    proportional to the degree of parallelism provided by the system.

5803                    It can vary from a few cycles when each competing thread is running on its own  
5804                    processor, to an indefinite amount of time when all threads are multiplexed on a single  
5805                    processor (which is why spin locking is not advisable on uniprocessors).

5806                    • When used properly, the amount of time the calling thread spends waiting on a spin lock  
5807                    should be considerably less than the time required to set up a corresponding watchdog  
5808                    timer. Since the primary purpose of spin locks is to provide a low-overhead  
5809                    synchronization mechanism for multi-processors, the overhead of a timeout mechanism  
5810                    was deemed unacceptable.

5811                    It was also suggested that an additional *count* argument be provided (on the  
5812                    *pthread\_spin\_lock()* call) in lieu of a true timeout so that a spin lock call could fail gracefully if  
5813                    it was unable to apply the lock after *count* attempts. This idea was rejected because it is not  
5814                    existing practice. Furthermore, the same effect can be obtained with *pthread\_spin\_trylock()*,  
5815                    as illustrated below:



```

5816         int n = MAX_SPIN;
5817         while ( --n >= 0 )
5818         {
5819             if ( !pthread_spin_try_lock(...) )
5820                 break;
5821         }
5822         if ( n >= 0 )
5823         {
5824             /* Successfully acquired the lock */
5825         }
5826         else
5827         {
5828             /* Unable to acquire the lock */
5829         }

```

5830 • *process-shared* Attribute

5831 The initialization functions associated with most POSIX synchronization objects (for  
5832 example, mutexes, barriers, and read-write locks) take an attributes object with a *process-*  
5833 *shared* attribute that specifies whether or not the object is to be shared across processes. In the  
5834 draft corresponding to the first balloting round, two separate initialization functions are  
5835 provided for spin locks, however: one for spin locks that were to be shared across processes  
5836 (*spin\_init()*), and one for locks that were only used by multiple threads within a single  
5837 process (*pthread\_spin\_init()*). This was done so as to keep the overhead associated with spin  
5838 waiting to an absolute minimum. However, the balloting group requested that, since the  
5839 overhead associated to a bit check was small, spin locks should be consistent with the rest of  
5840 the synchronization primitives, and thus the *process-shared* attribute was introduced for spin  
5841 locks.

5842 • Spin Locks versus Mutexes

5843 It has been suggested that mutexes are an adequate synchronization mechanism and spin  
5844 locks are not necessary. Locking mechanisms typically must trade off the processor resources  
5845 consumed while setting up to block the thread and the processor resources consumed by the  
5846 thread while it is blocked. Spin locks require very little resources to set up the blocking of a  
5847 thread. Existing practice is to simply loop, repeating the atomic locking operation until the  
5848 lock is available. While the resources consumed to set up blocking of the thread are low, the  
5849 thread continues to consume processor resources while it is waiting.

5850 On the other hand, mutexes may be implemented such that the processor resources  
5851 consumed to block the thread are large relative to a spin lock. After detecting that the mutex  
5852 lock is not available, the thread must alter its scheduling state, add itself to a set of waiting  
5853 threads, and, when the lock becomes available again, undo all of this before taking over  
5854 ownership of the mutex. However, while a thread is blocked by a mutex, no processor  
5855 resources are consumed.

5856 Therefore, spin locks and mutexes may be implemented to have different characteristics.  
5857 Spin locks may have lower overall overhead for very short-term blocking, and mutexes may  
5858 have lower overall overhead when a thread will be blocked for longer periods of time. The  
5859 presence of both interfaces allows implementations with these two different characteristics,  
5860 both of which may be useful to a particular application.

5861 It has also been suggested that applications can build their own spin locks from the  
5862 *pthread\_mutex\_trylock()* function:

5863           while (pthread\_mutex\_trylock(&mutex));

5864           The apparent simplicity of this construct is somewhat deceiving, however. While the actual  
5865           wait is quite efficient, various guarantees on the integrity of mutex objects (for example,  
5866           priority inheritance rules) may add overhead to the successful path of the trylock operation  
5867           that is not required of spin locks. One could, of course, add an attribute to the mutex to  
5868           bypass such overhead, but the very act of finding and testing this attribute represents more  
5869           overhead than is found in the typical spin lock.

5870           The need to hold spin lock overhead to an absolute minimum also makes it impossible to  
5871           provide guarantees against starvation similar to those provided for mutexes or read-write  
5872           locks. The overhead required to implement such guarantees (for example, disabling  
5873           preemption before spinning) may well exceed the overhead of the spin wait itself by many  
5874           orders of magnitude. If a “safe” spin wait seems desirable, it can always be provided (albeit  
5875           at some performance cost) via appropriate mutex attributes.

### 5876           **XSI Supported Functions**

5877           On XSI-conformant systems, the following symbolic constants are always defined:

5878            \_POSIX\_READER\_WRITER\_LOCKS  
5879            \_POSIX\_THREAD\_ATTR\_STACKADDR  
5880            \_POSIX\_THREAD\_ATTR\_STACKSIZE  
5881            \_POSIX\_THREAD\_PROCESS\_SHARED  
5882            \_POSIX\_THREADS

5883           Therefore, the following threads functions are always supported:

5884 <i>pthread_atfork()</i>	<i>pthread_key_delete()</i>
5885 <i>pthread_attr_destroy()</i>	<i>pthread_kill()</i>
5886 <i>pthread_attr_getdetachstate()</i>	<i>pthread_mutex_destroy()</i>
5887 <i>pthread_attr_getguardsize()</i>	<i>pthread_mutex_init()</i>
5888 <i>pthread_attr_getschedparam()</i>	<i>pthread_mutex_lock()</i>
5889 <i>pthread_attr_getstackaddr()</i>	<i>pthread_mutex_trylock()</i>
5890 <i>pthread_attr_getstacksize()</i>	<i>pthread_mutex_unlock()</i>
5891 <i>pthread_attr_init()</i>	<i>pthread_mutexattr_destroy()</i>
5892 <i>pthread_attr_setdetachstate()</i>	<i>pthread_mutexattr_getpshared()</i>
5893 <i>pthread_attr_setguardsize()</i>	<i>pthread_mutexattr_gettype()</i>
5894 <i>pthread_attr_setschedparam()</i>	<i>pthread_mutexattr_init()</i>
5895 <i>pthread_attr_setstackaddr()</i>	<i>pthread_mutexattr_setpshared()</i>
5896 <i>pthread_attr_setstacksize()</i>	<i>pthread_mutexattr_settype()</i>
5897 <i>pthread_cancel()</i>	<i>pthread_once()</i>
5898 <i>pthread_cleanup_pop()</i>	<i>pthread_rwlock_destroy()</i>
5899 <i>pthread_cleanup_push()</i>	<i>pthread_rwlock_init()</i>
5900 <i>pthread_cond_broadcast()</i>	<i>pthread_rwlock_rdlock()</i>
5901 <i>pthread_cond_destroy()</i>	<i>pthread_rwlock_tryrdlock()</i>
5902 <i>pthread_cond_init()</i>	<i>pthread_rwlock_trywrlock()</i>
5903 <i>pthread_cond_signal()</i>	<i>pthread_rwlock_unlock()</i>
5904 <i>pthread_cond_timedwait()</i>	<i>pthread_rwlock_wrlock()</i>
5905 <i>pthread_cond_wait()</i>	<i>pthread_rwlockattr_destroy()</i>
5906 <i>pthread_condattr_destroy()</i>	<i>pthread_rwlockattr_getpshared()</i>
5907 <i>pthread_condattr_getpshared()</i>	<i>pthread_rwlockattr_init()</i>
5908 <i>pthread_condattr_init()</i>	<i>pthread_rwlockattr_setpshared()</i>

5909	<i>pthread_condattr_setpshared()</i>	<i>pthread_self()</i>
5910	<i>pthread_create()</i>	<i>pthread_setcancelstate()</i>
5911	<i>pthread_detach()</i>	<i>pthread_setcanceltype()</i>
5912	<i>pthread_equal()</i>	<i>pthread_setconcurrency()</i>
5913	<i>pthread_exit()</i>	<i>pthread_setspecific()</i>
5914	<i>pthread_getconcurrency()</i>	<i>pthread_sigmask()</i>
5915	<i>pthread_getspecific()</i>	<i>pthread_testcancel()</i>
5916	<i>pthread_join()</i>	<i>sigwait()</i>
5917	<i>pthread_key_create()</i>	

5918 On XSI-conformant systems, the symbolic constant `_POSIX_THREAD_SAFE_FUNCTIONS` is  
5919 always defined. Therefore, the following functions are always supported:

5920	<i>asctime_r()</i>	<i>getpwnam_r()</i>
5921	<i>ctime_r()</i>	<i>getpwuid_r()</i>
5922	<i>flockfile()</i>	<i>gmtime_r()</i>
5923	<i>ftrylockfile()</i>	<i>localtime_r()</i>
5924	<i>funlockfile()</i>	<i>putc_unlocked()</i>
5925	<i>getc_unlocked()</i>	<i>putchar_unlocked()</i>
5926	<i>getchar_unlocked()</i>	<i>rand_r()</i>
5927	<i>getgrgid_r()</i>	<i>readdir_r()</i>
5928	<i>getgrnam_r()</i>	<i>strtok_r()</i>

5929 The following threads functions are only supported on XSI-conformant systems if the Realtime  
5930 Threads Option Group is supported :

5931	<i>pthread_attr_getinheritsched()</i>	<i>pthread_mutex_getprioceiling()</i>
5932	<i>pthread_attr_getschedpolicy()</i>	<i>pthread_mutex_setprioceiling()</i>
5933	<i>pthread_attr_getscope()</i>	<i>pthread_mutexattr_getprioceiling()</i>
5934	<i>pthread_attr_setinheritsched()</i>	<i>pthread_mutexattr_getprotocol()</i>
5935	<i>pthread_attr_setschedpolicy()</i>	<i>pthread_mutexattr_setprioceiling()</i>
5936	<i>pthread_attr_setscope()</i>	<i>pthread_mutexattr_setprotocol()</i>
5937	<i>pthread_getschedparam()</i>	<i>pthread_setschedparam()</i>

### 5938 XSI Threads Extensions

5939 The following XSI extensions to POSIX.1c are now supported in IEEE Std. 1003.1-200x as part of  
5940 the alignment with the Single UNIX Specification:

- 5941 • Extended mutex attribute types
- 5942 • Read-write locks and attributes (also introduced by IEEE Std. 1003.1j-2000 amendment)
- 5943 • Thread concurrency level
- 5944 • Thread stack guard size
- 5945 • Parallel I/O

5946 A total of 19 new functions were added.

5947 These extensions carefully follow the threads programming model specified in POSIX.1c. As  
5948 with POSIX.1c, all the new functions return zero if successful; otherwise, an error number is  
5949 returned to indicate the error.

5950 The concept of attribute objects was introduced in POSIX.1c to allow implementations to extend  
 5951 IEEE Std. 1003.1-200x without changing the existing interfaces. Attribute objects were defined  
 5952 for threads, mutexes, and condition variables. Attributes objects are defined as implementation-  
 5953 defined opaque types to aid extensibility, and functions are defined to allow attributes to be set  
 5954 or retrieved. This model has been followed when adding the new type attribute of  
 5955 **pthread\_mutexattr\_t** or the new read-write lock attributes object **pthread\_rwlockattr\_t**.

5956 • Extended Mutex Attributes

5957 POSIX.1c defines a mutex attributes object as an implementation-defined opaque object of  
 5958 type **pthread\_mutexattr\_t**, and specifies a number of attributes which this object must have  
 5959 and a number of functions which manipulate these attributes. These attributes include  
 5960 *detachstate*, *inheritsched*, *schedparm*, *schedpolicy*, *contentionscope*, *stackaddr*, and *stacksize*.

5961 The System Interfaces volume of IEEE Std. 1003.1-200x specifies another mutex attribute  
 5962 called *type*. The *type* attribute allows applications to specify the behavior of mutex locking  
 5963 operations in situations where the POSIX.1c behavior is undefined. The OSF DCE threads  
 5964 implementation, based on Draft 4 of POSIX.1c, specified a similar attribute. Note that the  
 5965 names of the attributes have changed somewhat from the OSF DCE threads implementation.

5966 The System Interfaces volume of IEEE Std. 1003.1-200x also extends the specification of the  
 5967 following POSIX.1c functions which manipulate mutexes:

5968 *pthread\_mutex\_lock()*  
 5969 *pthread\_mutex\_trylock()*  
 5970 *pthread\_mutex\_unlock()*

5971 to take account of the new mutex attribute type and to specify behavior which was declared  
 5972 as undefined in POSIX.1c. How a calling thread acquires or releases a mutex now depends  
 5973 upon the mutex *type* attribute.

5974 The *type* attribute can have the following values:

5975 PTHREAD\_MUTEX\_NORMAL

5976 Basic mutex with no specific error checking built in. Does not report a deadlock error.

5977 PTHREAD\_MUTEX\_RECURSIVE

5978 Allows any thread to recursively lock a mutex. The mutex must be unlocked an equal  
 5979 number of times to release the mutex.

5980 PTHREAD\_MUTEX\_ERRORCHECK

5981 Detects and reports simple usage errors; that is, an attempt to unlock a mutex that is not  
 5982 locked by the calling thread or that is not locked at all, or an attempt to relock a mutex  
 5983 the thread already owns.

5984 PTHREAD\_MUTEX\_DEFAULT

5985 The default mutex type. May be mapped to any of the above mutex types or may be an  
 5986 implementation-defined type.

5987 *Normal* mutexes do not detect deadlock conditions; for example, a thread will hang if it tries  
 5988 to relock a normal mutex that it already owns. Attempting to unlock a mutex locked by  
 5989 another thread, or unlocking an unlocked mutex, results in undefined behavior. Normal  
 5990 mutexes will usually be the fastest type of mutex available on a platform but provide the  
 5991 least error checking.

5992 *Recursive* mutexes are useful for converting old code where it is difficult to establish clear  
 5993 boundaries of synchronization. A thread can relock a recursive mutex without first unlocking  
 5994 it. The relocking deadlock which can occur with normal mutexes cannot occur with this type  
 5995 of mutex. However, multiple locks of a recursive mutex require the same number of unlocks

5996 to release the mutex before another thread can acquire the mutex. Furthermore, this type of  
5997 mutex maintains the concept of an owner. Thus, a thread attempting to unlock a recursive  
5998 mutex which another thread has locked returns with an error. A thread attempting to unlock  
5999 a recursive mutex that is not locked shall return with an error. Never use a recursive mutex  
6000 with condition variables because the implicit unlock performed by *pthread\_cond\_wait()* or  
6001 *pthread\_cond\_timedwait()* will not actually release the mutex if it had been locked multiple  
6002 times.

6003 *Errorcheck* mutexes provide error checking and are useful primarily as a debugging aid. A  
6004 thread attempting to relock an errorcheck mutex without first unlocking it returns with an  
6005 error. Again, this type of mutex maintains the concept of an owner. Thus, a thread  
6006 attempting to unlock an errorcheck mutex which another thread has locked returns with an  
6007 error. A thread attempting to unlock an errorcheck mutex that is not locked also returns with  
6008 an error. It should be noted that errorcheck mutexes will almost always be much slower than  
6009 normal mutexes due to the extra state checks performed.

6010 The *default* mutex type provides implementation-defined error checking. The default mutex  
6011 may be mapped to one of the other defined types or may be something entirely different.  
6012 This enables each vendor to provide the mutex semantics which the vendor feels will be  
6013 most useful to their target users. Most vendors will probably choose to make normal  
6014 mutexes the default so as to give applications the benefit of the fastest type of mutexes  
6015 available on their platform. Check your implementation's documentation.

6016 An application developer can use any of the mutex types almost interchangeably as long as  
6017 the application does not depend upon the implementation detecting (or failing to detect) any  
6018 particular errors. Note that a recursive mutex can be used with condition variable waits as  
6019 long as the application never recursively locks the mutex.

6020 Two functions are provided for manipulating the *type* attribute of a mutex attributes object.  
6021 This attribute is set or returned in the *type* parameter of these functions. The  
6022 *pthread\_mutexattr\_settype()* function is used to set a specific type value while  
6023 *pthread\_mutexattr\_gettype()* is used to return the type of the mutex. Setting the *type* attribute  
6024 of a mutex attributes object affects only mutexes initialized using that mutex attributes  
6025 object. Changing the *type* attribute does not affect mutexes previously initialized using that  
6026 mutex attributes object.

#### 6027 • Read-Write Locks and Attributes

6028 The read-write locks introduced have been harmonized with those in IEEE Std. 1003.1j-2000;  
6029 see also Section B.2.9.6 (on page 3472).

6030 Read-write locks (also known as reader-writer locks) allow a thread to exclusively lock some  
6031 shared data while updating that data, or allow any number of threads to have simultaneous  
6032 read-only access to the data.

6033 Unlike a mutex, a read-write lock distinguishes between reading data and writing data. A  
6034 mutex excludes all other threads. A read-write lock allows other threads access to the data,  
6035 providing no thread is modifying the data. Thus, a read-write lock is less primitive than  
6036 either a mutex-condition variable pair or a semaphore.

6037 Application developers should consider using a read-write lock rather than a mutex to  
6038 protect data that is frequently referenced but seldom modified. Most threads (readers) will be  
6039 able to read the data without waiting and will only have to block when some other thread (a  
6040 writer) is in the process of modifying the data. Conversely a thread that wants to change the  
6041 data is forced to wait until there are no readers. This type of lock is often used to facilitate  
6042 parallel access to data on multi-processor platforms or to avoid context switches on single  
6043 processor platforms where multiple threads access the same data.

6044 If a read-write lock becomes unlocked and there are multiple threads waiting to acquire the  
6045 write lock, the implementation's scheduling policy determines which thread shall acquire the  
6046 read-write lock for writing. If there are multiple threads blocked on a read-write lock for both  
6047 read locks and write locks, it is unspecified whether the readers or a writer acquire the lock  
6048 first. However, for performance reasons, implementations often favor writers over readers to  
6049 avoid potential writer starvation.

6050 A read-write lock object is an implementation-defined opaque object of type  
6051 **pthread\_rwlock\_t** as defined in `<pthread.h>`. There are two different sorts of locks  
6052 associated with a read-write lock: a *read lock* and a *write lock*.

6053 The `pthread_rwlockattr_init()` function initializes a read-write lock attributes object with the  
6054 default value for all the attributes defined in the implementation. After a read-write lock  
6055 attributes object has been used to initialize one or more read-write locks, changes to the  
6056 read-write lock attributes object, including destruction, do not affect previously initialized  
6057 read-write locks.

6058 Implementations must provide at least the read-write lock attribute *process-shared*. This  
6059 attribute can have the following values:

6060 PTHREAD\_PROCESS\_SHARED

6061 Any thread of any process that has access to the memory where the read-write lock  
6062 resides can manipulate the read-write lock.

6063 PTHREAD\_PROCESS\_PRIVATE

6064 Only threads created within the same process as the thread that initialized the read-  
6065 write lock can manipulate the read-write lock. This is the default value.

6066 The `pthread_rwlockattr_setpshared()` function is used to set the *process-shared* attribute of an  
6067 initialized read-write lock attributes object while the function `pthread_rwlockattr_getpshared()`  
6068 obtains the current value of the *process-shared* attribute.

6069 A read-write lock attributes object is destroyed using the `pthread_rwlockattr_destroy()`  
6070 function. The effect of subsequent use of the read-write lock attributes object is undefined.

6071 A thread creates a read-write lock using the `pthread_rwlock_init()` function. The attributes of  
6072 the read-write lock can be specified by the application developer; otherwise, the default  
6073 implementation-defined read-write lock attributes are used if the pointer to the read-write  
6074 lock attributes object is NULL. In cases where the default attributes are appropriate, the  
6075 PTHREAD\_RWLOCK\_INITIALIZER macro can be used to initialize statically allocated  
6076 read-write locks.

6077 A thread which wants to apply a read lock to the read-write lock can use either  
6078 `pthread_rwlock_rdlock()` or `pthread_rwlock_tryrdlock()`. If `pthread_rwlock_rdlock()` is used, the  
6079 thread acquires a read lock if a writer does not hold the write lock and there are no writers  
6080 blocked on the write lock. If a read lock is not acquired, the calling thread blocks until it can  
6081 acquire a lock. However, if `pthread_rwlock_tryrdlock()` is used, the function returns  
6082 immediately with the error [EBUSY] if any thread holds a write lock or there are blocked  
6083 writers waiting for the write lock.

6084 A thread which wants to apply a write lock to the read-write lock can use either of two  
6085 functions: `pthread_rwlock_wrlock()` or `pthread_rwlock_trywrlock()`. If `pthread_rwlock_wrlock()`  
6086 is used, the thread acquires the write lock if no other reader or writer threads hold the read-  
6087 write lock. If the write lock is not acquired, the thread blocks until it can acquire the write  
6088 lock. However, if `pthread_rwlock_trywrlock()` is used, the function returns immediately with  
6089 the error [EBUSY] if any thread is holding either a read or a write lock.

6090 The `pthread_rwlock_unlock()` function is used to unlock a read-write lock object held by the  
6091 calling thread. Results are undefined if the read-write lock is not held by the calling thread. If  
6092 there are other read locks currently held on the read-write lock object, the read-write lock  
6093 object shall remain in the read locked state but without the current thread as one of its  
6094 owners. If this function releases the last read lock for this read-write lock object, the read-  
6095 write lock object shall be put in the unlocked read state. If this function is called to release a  
6096 write lock for this read-write lock object, the read-write lock object shall be put in the  
6097 unlocked state.

#### 6098 • Thread Concurrency Level

6099 On threads implementations that multiplex user threads onto a smaller set of kernel  
6100 execution entities, the system attempts to create a reasonable number of kernel execution  
6101 entities for the application upon application startup.

6102 On some implementations, these kernel entities are retained by user threads that block in the  
6103 kernel. Other implementations do not *timeslice* user threads so that multiple compute-bound  
6104 user threads can share a kernel thread. On such implementations, some applications may use  
6105 up all the available kernel execution entities before its user-space threads are used up. The  
6106 process may be left with user threads capable of doing work for the application but with no  
6107 way to schedule them.

6108 The `pthread_setconcurrency()` function enables an application to request more kernel entities;  
6109 that is, specify a desired concurrency level. However, this function merely provides a hint to  
6110 the implementation. The implementation is free to ignore this request or to provide some  
6111 other number of kernel entities. If an implementation does not multiplex user threads onto a  
6112 smaller number of kernel execution entities, the `pthread_setconcurrency()` function has no  
6113 effect.

6114 The `pthread_setconcurrency()` function may also have an effect on implementations where the  
6115 kernel mode and user mode schedulers cooperate to ensure that ready user threads are not  
6116 prevented from running by other threads blocked in the kernel.

6117 The `pthread_getconcurrency()` function always returns the value set by a previous call to  
6118 `pthread_setconcurrency()`. However, if `pthread_setconcurrency()` was not previously called, this  
6119 function shall return zero to indicate that the threads implementation is maintaining the  
6120 concurrency level.

#### 6121 • Thread Stack Guard Size

6122 DCE threads introduced the concept of a *thread stack guard size*. Most thread  
6123 implementations add a region of protected memory to a thread's stack, commonly known as  
6124 a *guard region*, as a safety measure to prevent stack pointer overflow in one thread from  
6125 corrupting the contents of another thread's stack. The default size of the guard regions  
6126 attribute is {PAGESIZE} bytes and is implementation-defined.

6127 Some application developers may wish to change the stack guard size. When an application  
6128 creates a large number of threads, the extra page allocated for each stack may strain system  
6129 resources. In addition to the extra page of memory, the kernel's memory manager has to keep  
6130 track of the different protections on adjoining pages. When this is a problem, the application  
6131 developer may request a guard size of 0 bytes to conserve system resources by eliminating  
6132 stack overflow protection.

6133 Conversely an application that allocates large data structures such as arrays on the stack may  
6134 wish to increase the default guard size in order to detect stack overflow. If a thread allocates  
6135 two pages for a data array, a single guard page provides little protection against thread stack  
6136 overflows since the thread can corrupt adjoining memory beyond the guard page.

6137 The System Interfaces volume of IEEE Std. 1003.1-200x defines a new attribute of a thread  
6138 attributes object; that is, the *guardsize* attribute which allows applications to specify the size  
6139 of the guard region of a thread's stack.

6140 Two functions are provided for manipulating a thread's stack guard size. The  
6141 *pthread\_attr\_setguardsize()* function sets the thread *guardsize* attribute, and the  
6142 *pthread\_attr\_getguardsize()* function retrieves the current value.

6143 An implementation may round up the requested guard size to a multiple of the configurable  
6144 system variable {PAGESIZE}. In this case, *pthread\_attr\_getguardsize()* returns the guard size  
6145 specified by the previous *pthread\_attr\_setguardsize()* function call and not the rounded up  
6146 value.

6147 If an application is managing its own thread stacks using the *stackaddr* attribute, the *guardsize*  
6148 attribute is ignored and no stack overflow protection is provided. In this case, it is the  
6149 responsibility of the application to manage stack overflow along with stack allocation.

#### 6150 • Parallel I/O

6151 Many I/O intensive applications, such as database engines, attempt to improve performance  
6152 through the use of parallel I/O. However, POSIX.1 does not support parallel I/O very well  
6153 because the current offset of a file is an attribute of the file descriptor.

6154 Suppose two or more threads independently issue read requests on the same file. To read  
6155 specific data from a file, a thread must first call *lseek()* to seek to the proper offset in the file,  
6156 and then call *read()* to retrieve the required data. If more than one thread does this at the  
6157 same time, the first thread may complete its seek call, but before it gets a chance to issue its  
6158 read call a second thread may complete its seek call, resulting in the first thread accessing  
6159 incorrect data when it issues its read call. One workaround is to lock the file descriptor while  
6160 seeking and reading or writing, but this reduces parallelism and adds overhead.

6161 Instead, the System Interfaces volume of IEEE Std. 1003.1-200x provides two functions to  
6162 make seek/read and seek/write operations atomic. The file descriptor's current offset is  
6163 unchanged, thus allowing multiple read and write operations to proceed in parallel. This  
6164 improves the I/O performance of threaded applications. The *pread()* function is used to do  
6165 an atomic read of data from a file into a buffer. Conversely, the *pwrite()* function does an  
6166 atomic write of data from a buffer to a file.

#### 6167 B.2.9.1 Thread-Safety

6168 All functions required by IEEE Std. 1003.1-200x need to be thread-safe. Implementations have to  
6169 provide internal synchronization when necessary in order to achieve this goal. In certain  
6170 cases—for example, most floating-point implementations—context switch code may have to  
6171 manage the writable shared state.

6172 It is not required that all functions provided by IEEE Std. 1003.1-200x be either async-cancel-safe  
6173 or async-signal-safe.

6174 As it turns out, some functions are inherently not thread-safe; that is, their interface  
6175 specifications preclude reentrancy. For example, some functions (such as *asctime()*) return a  
6176 pointer to a result stored in memory space allocated by the function on a per-process basis. Such  
6177 a function is not thread-safe, because its result can be overwritten by successive invocations.  
6178 Other functions, while not inherently non-thread-safe, may be implemented in ways that lead to  
6179 them not being thread-safe. For example, some functions (such as *rand()*) store state information  
6180 (such as a seed value, which survives multiple function invocations) in memory space allocated  
6181 by the function on a per-process basis. The implementation of such a function is not thread-safe  
6182 if the implementation fails to synchronize invocations of the function and thus fails to protect



6183 the state information. The problem is that when the state information is not protected,  
6184 concurrent invocations can interfere with one another (for example, see the same seed value).

#### 6185 *Thread-Safety and Locking of Existing Functions*

6186 Originally, POSIX.1 was not designed to work in a multi-threaded environment, and some  
6187 implementations of some existing functions will not work properly when executed concurrently.  
6188 To provide routines that will work correctly in an environment with threads (“thread-safe”), two  
6189 problems need to be solved:

- 6190 1. Routines that maintain or return pointers to static areas internal to the routine (which may  
6191 now be shared) need to be modified. The routines *ttyname()* and *localtime()* are examples.
- 6192 2. Routines that access data space shared by more than one thread need to be modified. The  
6193 *malloc()* function and the *stdio* family routines are examples.

6194 There are a variety of constraints on these changes. The first is compatibility with the existing  
6195 versions of these functions—non-thread-safe functions will continue to be in use for some time,  
6196 as the original interfaces are used by existing code. Another is that the new thread-safe versions  
6197 of these functions represent as small a change as possible over the familiar interfaces provided  
6198 by the existing non-thread-safe versions. The new interfaces should be independent of any  
6199 particular threads implementation. In particular, they should be thread-safe without depending  
6200 on explicit thread-specific memory. Finally, there should be minimal performance penalty due to  
6201 the changes made to the functions.

6202 It is intended that the list of functions from POSIX.1 that cannot be made thread-safe and for  
6203 which corrected versions are provided be complete.

#### 6204 *Thread-Safety and Locking Solutions*

6205 Many of the POSIX.1 functions were thread-safe and did not change at all. However, some  
6206 functions (for example, the math functions typically found in **libm**) are not thread-safe because  
6207 of writable shared global state. For instance, in IEEE Std. 754-1985 floating-point  
6208 implementations, the computation modes and flags are global and shared.

6209 Some functions are not thread-safe because a particular implementation is not reentrant,  
6210 typically because of a non-essential use of static storage. These require only a new  
6211 implementation.

6212 Thread-safe libraries are useful in a wide range of parallel (and asynchronous) programming  
6213 environments, not just within pthreads. In order to be used outside the context of pthreads,  
6214 however, such libraries still have to use some synchronization method. These could either be  
6215 independent of the pthread synchronization operations, or they could be a subset of the pthread  
6216 interfaces. Either method results in thread-safe library implementations that can be used without  
6217 the rest of pthreads.

6218 Some functions, such as the *stdio* family interface and dynamic memory allocation functions  
6219 such as *malloc()*, are interdependent routines that share resources (for example, buffers) across  
6220 related calls. These require synchronization to work correctly, but they do not require any  
6221 change to their external (user-visible) interfaces.

6222 In some cases, such as *getc()* and *putc()*, adding synchronization is likely to create an  
6223 unacceptable performance impact. In this case, slower thread-safe synchronized functions are to  
6224 be provided, but the original, faster (but unsafe) functions (which may be implemented as  
6225 macros) are retained under new names. Some additional special-purpose synchronization  
6226 facilities are necessary for these macros to be usable in multi-threaded programs. This also  
6227 requires changes in `<stdio.h>`.

6228 The other common reason that functions are unsafe is that they return a pointer to static storage,  
6229 making the functions non-thread-safe. This has to be changed, and there are three natural  
6230 choices:

6231 1. Return a pointer to thread-specific storage

6232 This could incur a severe performance penalty on those architectures with a costly  
6233 implementation of the thread-specific data interface.

6234 A variation on this technique is to use *malloc()* to allocate storage for the function output  
6235 and return a pointer to this storage. This technique may also have an undesirable  
6236 performance impact, however, and a simplistic implementation requires that the user  
6237 program explicitly free the storage object when it is no longer needed. This technique is  
6238 used by some existing POSIX.1 functions. With careful implementation for infrequently  
6239 used functions, there may be little or no performance or storage penalty, and the  
6240 maintenance of already-standardized interfaces is a significant benefit.

6241 2. Return the actual value computed by the function

6242 This technique can only be used with functions that return pointers to structures—routines  
6243 that return character strings would have to wrap their output in an enclosing structure in  
6244 order to return the output on the stack. There is also a negative performance impact  
6245 inherent in this solution in that the output value has to be copied twice before it can be  
6246 used by the calling function: once from the called routine's local buffers to the top of the  
6247 stack, then from the top of the stack to the assignment target. Finally, many older  
6248 compilers cannot support this technique due to a historical tendency to use internal static  
6249 buffers to deliver the results of structure-valued functions.

6250 3. Have the caller pass the address of a buffer to contain the computed value

6251 The only disadvantage of this approach is that extra arguments have to be provided by the  
6252 calling program. It represents the most efficient solution to the problem, however, and,  
6253 unlike the *malloc()* technique, it is semantically clear.

6254 There are some routines (often groups of related routines) whose interfaces are inherently non-  
6255 thread-safe because they communicate across multiple function invocations by means of static  
6256 memory locations. The solution is to redesign the calls so that they are thread-safe, typically by  
6257 passing the needed data as extra parameters. Unfortunately, this may require major changes to  
6258 the interface as well.

6259 A floating-point implementation using IEEE Std. 754-1985 is a case in point. A less problematic  
6260 example is the *rand48* family of pseudo-random number generators. The functions *getgrgid()*,  
6261 *getgrnam()*, *getpwnam()*, and *getpwuid()* are another such case.

6262 The problems with *errno* are discussed in **Alternative Solutions for Per-Thread *errno*** (on page  
6263 3390).

6264 Some functions can be thread-safe or not, depending on their arguments. These include the  
6265 *tmpnam()* and *ctermid()* functions. These functions have pointers to character strings as  
6266 arguments. If the pointers are not NULL, the functions store their results in the character string;  
6267 however, if the pointers are NULL, the functions store their results in an area that may be static  
6268 and thus subject to overwriting by successive calls. These should only be called by multi-thread  
6269 applications when their arguments are non-NULL.

6270 *Asynchronous Safety and Thread-Safety*

6271 A floating-point implementation has many modes that effect rounding and other aspects of  
6272 computation. Functions in some math library implementations may change the computation  
6273 modes for the duration of a function call. If such a function call is interrupted by a signal or

- 6274 cancelation, the floating-point state is not required to be protected.
- 6275 There is a significant cost to make floating-point operations async-cancel-safe or async-signal-  
6276 safe; accordingly, neither form of async safety is required.
- 6277 *Functions Returning Pointers to Static Storage*
- 6278 For those functions that are not thread-safe because they return values in fixed size statically  
6279 allocated structures, alternate “\_r” forms are provided that pass a pointer to an explicit result  
6280 structure. Those that return pointers into library-allocated buffers have forms provided with  
6281 explicit buffer and length parameters.
- 6282 For functions that return pointers to library-allocated buffers, it makes sense to provide “\_r”  
6283 versions that allow the application control over allocation of the storage in which results are  
6284 returned. This allows the state used by these functions to be managed on an application-specific  
6285 basis, supporting per-thread, per-process, or other application-specific sharing relationships.
- 6286 Early proposals had provided “\_r” versions for functions that returned pointers to variable-size  
6287 buffers without providing a means for determining the required buffer size. This would have  
6288 made using such functions exceedingly clumsy, potentially requiring iteratively calling them  
6289 with increasingly larger guesses for the amount of storage required. Hence, *sysconf()* variables  
6290 have been provided for such functions that return the maximum required buffer size.
- 6291 Thus, the rule that has been followed by IEEE Std. 1003.1-200x when adapting single-threaded  
6292 non-thread-safe library functions is as follows: all functions returning pointers to library-  
6293 allocated storage should have “\_r” versions provided, allowing the application control over the  
6294 storage allocation. Those with variable-sized return values accept both a buffer address and a  
6295 length parameter. The *sysconf()* variables are provided to supply the appropriate buffer sizes  
6296 when required. Implementors are encouraged to apply the same rule when adapting their own  
6297 existing functions to a pthreads environment.
- 6298 **B.2.9.2 Thread IDs**
- 6299 Separate programs should communicate through well-defined interfaces and should not depend  
6300 on each other's implementation. For example, if a programmer decides to rewrite the *sort*  
6301 program using multiple threads, it should be easy to do this so that the interface to the *sort*  
6302 program does not change. Consider that if the user causes SIGINT to be generated while the *sort*  
6303 program is running, keeping the same interface means that the entire sort program is killed, not  
6304 just one of its threads. As another example, consider a realtime program that manages a reactor.  
6305 Such a program may wish to allow other programs to control the priority at which it watches the  
6306 control rods. One technique to accomplish this is to write the ID of the thread watching the  
6307 control rods into a file and allow other programs to change the priority of that thread as they see  
6308 fit. A simpler technique is to have the reactor process accept IPCs (Inter-Process Communication  
6309 messages) from other processes, telling it at a semantic level what priority the program should  
6310 assign to watching the control rods. This allows the programmer greater flexibility in the  
6311 implementation. For example, the programmer can change the implementation from having one  
6312 thread per rod to having one thread watching all of the rods without changing the interface.  
6313 Having threads live inside the process means that the implementation of a process is invisible to  
6314 outside processes (excepting debuggers and system management tools).
- 6315 Threads do not provide a protection boundary. Every thread model allows threads to share  
6316 memory with other threads and encourages this sharing to be widespread. This means that one  
6317 thread can wipe out memory that is needed for the correct functioning of other threads that are  
6318 sharing its memory. Consequently, providing each thread with its own user and/or group IDs  
6319 would not provide a protection boundary between threads sharing memory.

6320 *B.2.9.3 Thread Mutexes*

6321 There is no additional rationale for this section.

6322 *B.2.9.4 Thread Scheduling*

## 6323 • Scheduling Implementation Models

6324 The following scheduling implementation models are presented in terms of threads and  
6325 “kernel entities”. This is to simplify exposition of the models, and it does not imply that an  
6326 implementation actually has an identifiable “kernel entity”.

6327 A kernel entity is not defined beyond the fact that it has scheduling attributes that are used to  
6328 resolve contention with other kernel entities for execution resources. A kernel entity may be  
6329 thought of as an envelope that holds a thread or a separate kernel thread. It is not a  
6330 conventional process, although it shares with the process the attribute that it has a single  
6331 thread of control; it does not necessarily imply an address space, open files, and so on. It is  
6332 better thought of as a primitive facility upon which conventional processes and threads may  
6333 be constructed.

## 6334 — System Thread Scheduling Model

6335 This model consists of one thread per kernel entity. The kernel entity is solely responsible  
6336 for scheduling thread execution on one or more processors. This model schedules all  
6337 threads against all other threads in the system using the scheduling attributes of the  
6338 thread.

## 6339 — Process Scheduling Model

6340 A generalized process scheduling model consists of two levels of scheduling. A threads  
6341 library creates a pool of kernel entities, as required, and schedules threads to run on them  
6342 using the scheduling attributes of the threads. Typically, the size of the pool is a function  
6343 of the simultaneously runnable threads, not the total number of threads. The kernel then  
6344 schedules the kernel entities onto processors according to their scheduling attributes,  
6345 which are managed by the threads library. This set model potentially allows a wide range  
6346 of mappings between threads and kernel entities.

## 6347 • System and Process Scheduling Model Performance

6348 There are a number of important implications on the performance of applications using these  
6349 scheduling models. The process scheduling model potentially provides lower overhead for  
6350 making scheduling decisions, since there is no need to access kernel-level information or  
6351 functions and the set of schedulable entities is smaller (only the threads within the process).

6352 On the other hand, since the kernel is also making scheduling decisions regarding the system  
6353 resources under its control (for example, CPU(s), I/O devices, memory), decisions that do  
6354 not take thread scheduling parameters into account can result in indeterminate delays for  
6355 realtime application threads, causing them to miss maximum response time limits.

## 6356 • Rate Monotonic Scheduling

6357 Rate monotonic scheduling was considered, but rejected for standardization in the context of  
6358 pthreads. A sporadic server policy is included.

## 6359 • Scheduling Options

6360 In IEEE Std. 1003.1-200x, the basic thread scheduling functions are defined under the Threads  
6361 option, so that they are required of all threads implementations. However, there are no  
6362 specific scheduling policies required by this option to allow for conforming thread  
6363 implementations that are not targeted to realtime applications.

6364 Specific standard scheduling policies are defined to be under the Thread Execution  
6365 Scheduling option, and they are specifically designed to support realtime applications by  
6366 providing predictable resource sharing sequences. The name of this option was chosen to  
6367 emphasize that this functionality is defined as appropriate for realtime applications that  
6368 require simple priority-based scheduling.

6369 It is recognized that these policies are not necessarily satisfactory for some multi-processor  
6370 implementations, and work is ongoing to address a wider range of scheduling behaviors. The  
6371 interfaces have been chosen to create abundant opportunity for future scheduling policies to  
6372 be implemented and standardized based on this interface. In order to standardize a new  
6373 scheduling policy, all that is required (from the standpoint of thread scheduling attributes) is  
6374 to define a new policy name, new members of the thread attributes object, and functions to  
6375 set these members when the scheduling policy is equal to the new value.

### 6376 **Scheduling Contention Scope**

6377 In order to accommodate the requirement for realtime response, each thread has a scheduling  
6378 contention scope attribute. Threads with a system scheduling contention scope have to be  
6379 scheduled with respect to all other threads in the system. These threads are usually bound to a  
6380 single kernel entity that reflects their scheduling attributes and are directly scheduled by the  
6381 kernel.

6382 Threads with a process scheduling contention scope need be scheduled only with respect to the  
6383 other threads in the process. These threads may be scheduled within the process onto a pool of  
6384 kernel entities. The implementation is also free to bind these threads directly to kernel entities  
6385 and let them be scheduled by the kernel. Process scheduling contention scope allows the  
6386 implementation the most flexibility and is the default if both contention scopes are supported  
6387 and none is specified.

6388 Thus, the choice by implementors to provide one or the other (or both) of these scheduling  
6389 models is driven by the need of their supported application domains for worst-case (that is,  
6390 realtime) response, or average-case (non-realtime) response.

### 6391 **Scheduling Allocation Domain**

6392 The SCHED\_FIFO and SCHED\_RR scheduling policies take on different characteristics on a  
6393 multi-processor. Other scheduling policies are also subject to changed behavior when executed  
6394 on a multi-processor. The concept of scheduling allocation domain determines the set of  
6395 processors on which the threads of an application may run. By considering the application's  
6396 processor scheduling allocation domain for its threads, scheduling policies can be defined in  
6397 terms of their behavior for varying processor scheduling allocation domain values. It is  
6398 conceivable that not all scheduling allocation domain sizes make sense for all scheduling  
6399 policies on all implementations. The concept of scheduling allocation domain, however, is a  
6400 useful tool for the description of multi-processor scheduling policies.

6401 The "process control" approach to scheduling obtains significant performance advantages from  
6402 dynamic scheduling allocation domain sizes when it is applicable.

6403 Non-Uniform Memory Access (NUMA) multi-processors may use a system scheduling structure  
6404 that involves reassignment of threads among scheduling allocation domains. In NUMA  
6405 machines, a natural model of scheduling is to match scheduling allocation domains to clusters of  
6406 processors. Load balancing in such an environment requires changing the scheduling allocation  
6407 domain to which a thread is assigned.

**6408 Scheduling Documentation**

6409 Implementation-provided scheduling policies need to be completely documented in order to be  
6410 useful. This documentation includes a description of the attributes required for the policy, the  
6411 scheduling interaction of threads running under this policy and all other supported policies, and  
6412 the effects of all possible values for processor scheduling allocation domain. Note that for the  
6413 implementor wishing to be minimally-compliant, it is (minimally) acceptable to define the  
6414 behavior as undefined.

**6415 Scheduling Contention Scope Attribute**

6416 The scheduling contention scope defines how threads compete for resources. Within  
6417 IEEE Std. 1003.1-200x, scheduling contention scope is used to describe only how threads are  
6418 scheduled in relation to one another in the system. That is, either they are scheduled against all  
6419 other threads in the system (“system scope”) or only against those threads in the process  
6420 (“process scope”). In fact, scheduling contention scope may apply to additional resources,  
6421 including virtual timers and profiling, which are not currently considered by  
6422 IEEE Std. 1003.1-200x.

**6423 Mixed Scopes**

6424 If only one scheduling contention scope is supported, the scheduling decision is straightforward.  
6425 To perform the processor scheduling decision in a mixed scope environment, it is necessary to  
6426 map the scheduling attributes of the thread with process-wide contention scope to the same  
6427 attribute space as the thread with system-wide contention scope.

6428 Since a conforming implementation has to support one and may support both scopes, it is useful  
6429 to discuss the effects of such choices with respect to example applications. If an implementation  
6430 supports both scopes, mixing scopes provides a means of better managing system-level (that is,  
6431 kernel-level) and library-level resources. In general, threads with system scope will require the  
6432 resources of a separate kernel entity in order to guarantee the scheduling semantics. On the  
6433 other hand, threads with process scope can share the resources of a kernel entity while  
6434 maintaining the scheduling semantics.

6435 The application is free to create threads with dedicated kernel resources, and other threads that  
6436 multiplex kernel resources. Consider the example of a window server. The server allocates two  
6437 threads per widget: one thread manages the widget user interface (including drawing), while the  
6438 other thread takes any required application action. This allows the widget to be “active” while  
6439 the application is computing. A screen image may be built from thousands of widgets. If each of  
6440 these threads had been created with system scope, then most of the kernel-level resources might  
6441 be wasted, since only a few widgets are active at any one time. In addition, mixed scope is  
6442 particularly useful in a window server where one thread with high priority and system scope  
6443 handles the mouse so that it tracks well. As another example, consider a database server. For  
6444 each of the hundreds or thousands of clients supported by a large server, an equivalent number  
6445 of threads will have to be created. If each of these threads were system, the consequences would  
6446 be the same as for the window server example above. However, the server could be constructed  
6447 so that actual retrieval of data is done by several dedicated threads. Dedicated threads that do  
6448 work for all clients frequently justify the added expense of system scope. If it were not  
6449 permissible to mix system and process threads in the same process, this type of solution would  
6450 not be possible.

**6451 Dynamic Thread Scheduling Parameters Access**

6452 In many time-constrained applications, there is no need to change the scheduling attributes  
6453 dynamically during thread or process execution, since the general use of these attributes is to  
6454 reflect directly the time constraints of the application. Since these time constraints are generally  
6455 imposed to meet higher-level system requirements, such as accuracy or availability, they  
6456 frequently should remain unchanged during application execution.

6457 However, there are important situations in which the scheduling attributes should be changed.  
6458 Generally, this will occur when external environmental conditions exist in which the time  
6459 constraints change. Consider, for example, a space vehicle major mode change, such as the  
6460 change from ascent to descent mode, or the change from the space environment to the  
6461 atmospheric environment. In such cases, the frequency with which many of the sensors or  
6462 acutators need to be read or written will change, which will necessitate a priority change. In  
6463 other cases, even the existence of a time constraint might be temporary, necessitating not just a  
6464 priority change, but also a policy change for ongoing threads or processes. For this reason, it is  
6465 critical that the interface should provide functions to change the scheduling parameters  
6466 dynamically, but, as with many of the other realtime functions, it is important that applications  
6467 use them properly to avoid the possibility of unnecessarily degrading performance.

6468 In providing functions for dynamically changing the scheduling behavior of threads, there were  
6469 two options: provide functions to get and set the individual scheduling parameters of threads, or  
6470 provide a single interface to get and set all the scheduling parameters for a given thread  
6471 simultaneously. Both approaches have merit. Access functions for individual parameters allow  
6472 simpler control of thread scheduling for simple thread scheduling parameters. However, a single  
6473 function for setting all the parameters for a given scheduling policy is required when first setting  
6474 that scheduling policy. Since the single all-encompassing functions are required, it was decided  
6475 to leave the interface as minimal as possible. Note that simpler functions (such as  
6476 *pthread\_setprio()* for threads running under the priority-based schedulers) can be easily defined  
6477 in terms of the all-encompassing functions.

6478 If the *pthread\_setschedparam()* function executes successfully, it will have set all of the scheduling  
6479 parameter values indicated in *param*; otherwise, none of the scheduling parameters will have  
6480 been modified. This is necessary to ensure that the scheduling of this and all other threads  
6481 continues to be consistent in the presence of an erroneous scheduling parameter.

6482 The [EPERM] error value is included in the list of possible *pthread\_setschedparam()* error returns  
6483 as a reflection of the fact that the ability to change scheduling parameters increases risks to the  
6484 implementation and application performance if the scheduling parameters are changed  
6485 improperly. For this reason, and based on some existing practice, it was felt that some  
6486 implementations would probably choose to define specific permissions for changing either a  
6487 thread's own or another thread's scheduling parameters. IEEE Std. 1003.1-200x does not include  
6488 portable methods for setting or retrieving permissions, so any such use of permissions is  
6489 completely unspecified .

**6490 Mutex Initialization Scheduling Attributes**

6491 In a priority-driven environment, a direct use of traditional primitives like mutexes and  
6492 condition variables can lead to unbounded priority inversion, where a higher priority thread can  
6493 be blocked by a lower priority thread, or set of threads, for an unbounded duration of time. As a  
6494 result, it becomes impossible to guarantee thread deadlines. Priority inversion can be bounded  
6495 and minimized by the use of priority inheritance protocols. This allows thread deadlines to be  
6496 guaranteed even in the presence of synchronization requirements.

6497 Two useful but simple members of the family of priority inheritance protocols are the basic  
6498 priority inheritance protocol and the priority ceiling protocol emulation. Under the Basic Priority

6499 Inheritance protocol (governed by the Threads Priority Inheritance option), a thread that is  
6500 blocking higher priority threads executes at the priority of the highest priority thread that it  
6501 blocks. This simple mechanism allows priority inversion to be bounded by the duration of  
6502 critical sections and makes timing analysis possible.

6503 Under the Priority Ceiling Protocol Emulation protocol (governed by the Thread Priority  
6504 Protection option), each mutex has a priority ceiling, usually defined as the priority of the  
6505 highest priority thread that can lock the mutex. When a thread is executing inside critical  
6506 sections, its priority is unconditionally increased to the highest of the priority ceilings of all the  
6507 mutexes owned by the thread. This protocol has two very desirable properties in uni-processor  
6508 systems. First, a thread can be blocked by a lower priority thread for at most the duration of one  
6509 single critical section. Furthermore, when the protocol is correctly used in a single processor, and  
6510 if threads do not become blocked while owning mutexes, mutual deadlocks are prevented.

6511 The priority ceiling emulation can be extended to multiple processor environments, in which  
6512 case the values of the priority ceilings will be assigned depending on the kind of mutex that is  
6513 being used: local to only one processor, or global, shared by several processors. Local priority  
6514 ceilings will be assigned the usual way, equal to the priority of the highest priority thread that  
6515 may lock that mutex. Global priority ceilings will usually be assigned a priority level higher than  
6516 all the priorities assigned to any of the threads that reside in the involved processors to avoid the  
6517 effect called remote blocking.

#### 6518 **Change the Priority Ceiling of a Mutex**

6519 In order for the priority protect protocol to exhibit its desired properties of bounding priority  
6520 inversion and avoidance of deadlock, it is critical that the ceiling priority of a mutex be the same  
6521 as the priority of the highest thread that can ever hold it, or higher. Thus, if the priorities of the  
6522 threads using such mutexes never change dynamically, there is no need ever to change the  
6523 priority ceiling of a mutex.

6524 However, if a major system mode change results in an altered response time requirement for one  
6525 or more application threads, their priority has to change to reflect it. It will occasionally be the  
6526 case that the priority ceilings of mutexes held also need to change. While changing priority  
6527 ceilings should generally be avoided, it is important that IEEE Std. 1003.1-200x provide these  
6528 interfaces for those cases in which it is necessary.

#### 6529 *B.2.9.5 Thread Cancellation*

6530 Many existing threads packages have facilities for canceling an operation or canceling a thread.  
6531 These facilities are used for implementing user requests (such as the CANCEL button in a  
6532 window-based application), for implementing OR parallelism (for example, telling the other  
6533 threads to stop working once one thread has found a forced mate in a parallel chess program), or  
6534 for implementing the ABORT mechanism in Ada.

6535 POSIX programs traditionally have used the signal mechanism combined with either *longjmp()*  
6536 or polling to cancel operations. Many POSIX programmers have trouble using these facilities to  
6537 solve their problems efficiently in a single-threaded process. With the introduction of threads,  
6538 these solutions become even more difficult to use.

6539 The main issues with implementing a cancellation facility are specifying the operation to be  
6540 canceled, cleanly releasing any resources allocated to that operation, controlling when the target  
6541 notices that it has been canceled, and defining the interaction between asynchronous signals and  
6542 cancellation.



**6543 Specifying the Operation to Cancel**

6544 Consider a thread that calls through five distinct levels of program abstraction and then, inside  
6545 the lowest-level abstraction, calls a function that suspends the thread. (An abstraction boundary  
6546 is a layer at which the client of the abstraction sees only the service being provided and can  
6547 remain ignorant of the implementation. Abstractions are often layered, each level of abstraction  
6548 being a client of the lower-level abstraction and implementing a higher-level abstraction.)  
6549 Depending on the semantics of each abstraction, one could imagine wanting to cancel only the  
6550 call that causes suspension, only the bottom two levels, or the operation being done by the entire  
6551 thread. Canceling operations at a finer grain than the entire thread is difficult because threads  
6552 are active and they may be run in parallel on a multi-processor. By the time one thread can make  
6553 a request to cancel an operation, the thread performing the operation may have completed that  
6554 operation and gone on to start another operation whose cancelation is not desired. Thread IDs  
6555 are not reused until the thread has exited, and either it was created with the *Attr detachstate*  
6556 attribute set to *PTHREAD\_CREATE\_DETACHED* or the *pthread\_join()* or *pthread\_detach()*  
6557 function has been called for that thread. Consequently, a thread cancelation will never be  
6558 misdirected when the thread terminates. For these reasons, the canceling of operations is done at  
6559 the granularity of the thread. Threads are designed to be inexpensive enough so that a separate  
6560 thread may be created to perform each separately cancelable operation; for example, each  
6561 possibly long running user request.

6562 For cancelation to be used in existing code, cancelation scopes and handlers will have to be  
6563 established for code that needs to release resources upon cancelation, so that it follows the  
6564 programming discipline described in the text.

**6565 A Special Signal Versus a Special Interface**

6566 Two different mechanisms were considered for providing the cancelation interfaces. The first  
6567 was to provide an interface to direct signals at a thread and then to define a special signal that  
6568 had the required semantics. The other alternative was to use a special interface that delivered the  
6569 correct semantics to the target thread.

6570 The solution using signals produced a number of problems. It required the implementation to  
6571 provide cancelation in terms of signals whereas a perfectly valid (and possibly more efficient)  
6572 implementation could have both layered on a low-level set of primitives. There were so many  
6573 exceptions to the special signal (it cannot be used with kill, no POSIX.1 interfaces can be used  
6574 with it) that it was clearly not a valid signal. Its semantics on delivery were also completely  
6575 different from any existing POSIX.1 signal. As such, a special interface that did not mandate the  
6576 implementation and did not confuse the semantics of signals and cancelation was felt to be the  
6577 better solution.

**6578 Races Between Cancelation and Resuming Execution**

6579 Due to the nature of cancelation, there is generally no synchronization between the thread  
6580 requesting the cancelation of a blocked thread and events that may cause that thread to resume  
6581 execution. For this reason, and because excess serialization hurts performance, when both an  
6582 event that a thread is waiting for has occurred and a cancelation request has been made and  
6583 cancelation is enabled, IEEE Std. 1003.1-200x explicitly allows the implementation to choose  
6584 between returning from the blocking call or acting on the cancelation request.

6585 **Interaction of Cancellation with Asynchronous Signals**

6586 A typical use of cancellation is to acquire a lock on some resource and to establish a cancellation  
6587 cleanup handler for releasing the resource when and if the thread is canceled.

6588 A correct and complete implementation of cancellation in the presence of asynchronous signals  
6589 requires considerable care. An implementation has to push a cancellation cleanup handler on the  
6590 cancellation cleanup stack while maintaining the integrity of the stack data structure. If an  
6591 asynchronously generated signal is posted to the thread during a stack operation, the signal  
6592 handler cannot manipulate the cancellation cleanup stack. As a consequence, asynchronous  
6593 signal handlers may not cancel threads or otherwise manipulate the cancellation state of a thread.  
6594 Threads may, of course, be canceled by another thread that used a *sigwait()* function to wait  
6595 synchronously for an asynchronous signal.

6596 In order for cancellation to function correctly, it is required that asynchronous signal handlers not  
6597 change the cancellation state. This requires that some elements of existing practice, such as using  
6598 *longjmp()* to exit from an asynchronous signal handler implicitly, be prohibited in cases where  
6599 the integrity of the cancellation state of the interrupt thread cannot be ensured.

6600 **Thread Cancellation Overview**

## 6601 • Cancellability States

6602 The three possible cancellability states (disabled, deferred, and asynchronous) are encoded  
6603 into two separate bits ((disable, enable) and (deferred, asynchronous)) to allow them to be  
6604 changed and restored independently. For instance, short code sequences that will not block  
6605 sometimes disable cancellability on entry and restore the previous state upon exit. Likewise,  
6606 long or unbounded code sequences containing no convenient explicit cancellation points will  
6607 sometimes set the cancellability type to asynchronous on entry and restore the previous value  
6608 upon exit.

## 6609 • Cancellation Points

6610 Cancellation points are points inside of certain functions where a thread has to act on any  
6611 pending cancellation request when cancellability is enabled, if the function would block. As  
6612 with checking for signals, operations need only check for pending cancellation requests when  
6613 the operation is about to block indefinitely.

6614 The idea was considered of allowing implementations to define whether blocking calls such  
6615 as *read()* should be cancellation points. It was decided that it would adversely affect the  
6616 design of portable applications if blocking calls were not cancellation points because threads  
6617 could be left blocked in an uncancelable state.

6618 There are several important blocking routines that are specifically not made cancellation  
6619 points:

6620 — *pthread\_mutex\_lock()*

6621 If *pthread\_mutex\_lock()* were a cancellation point, every routine that called it would also  
6622 become a cancellation point (that is, any routine that touched shared state would  
6623 automatically become a cancellation point). For example, *malloc()*, *free()*, and *rand()*  
6624 would become cancellation points under this scheme. Having too many cancellation points  
6625 makes programming very difficult, leading to either much disabling and restoring of  
6626 cancellability or much difficulty in trying to arrange for reliable cleanup at every possible  
6627 place.

6628 Since *pthread\_mutex\_lock()* is not a cancellation point, threads could result in being  
6629 blocked uninterruptibly for long periods of time if mutexes were used as a general

6630 synchronization mechanism. As this is normally not acceptable, mutexes should only be  
 6631 used to protect resources that are held for small fixed lengths of time where not being  
 6632 able to be canceled will not be a problem. Resources that need to be held exclusively for  
 6633 long periods of time should be protected with condition variables.

6634 — *barrier\_wait()*

6635 Canceling a barrier wait will render a barrier unusable. Similar to a barrier timeout (which  
 6636 the standard developers rejected), there is no way to guarantee the consistency of a  
 6637 barrier's internal data structures if a barrier wait is canceled.

6638 — *pthread\_spin\_lock()*

6639 As with mutexes, spin locks should only be used to protect resources that are held for  
 6640 small fixed lengths of time where not being cancelable will not be a problem.

6641 Every library routine should specify whether or not it includes any cancellation points.  
 6642 Typically, only those routines that may block or compute indefinitely need to include  
 6643 cancellation points.

6644 Correctly coded routines only reach cancellation points after having set up a cancellation  
 6645 cleanup handler to restore invariants if the thread is canceled at that point. Being cancelable  
 6646 only at specified cancellation points allows programmers to keep track of actions needed in a  
 6647 cancellation cleanup handler more easily. A thread should only be made asynchronously  
 6648 cancelable when it is not in the process of acquiring or releasing resources or otherwise in a  
 6649 state from which it would be difficult or impossible to recover.

#### 6650 • Thread Cancellation Cleanup Handlers

6651 The cancellation cleanup handlers provide a portable mechanism, easy to implement, for  
 6652 releasing resources and restoring invariants. They are easier to use than signal handlers  
 6653 because they provide a stack of cancellation cleanup handlers rather than a single handler,  
 6654 and because they have an argument that can be used to pass context information to the  
 6655 handler.

6656 The alternative to providing these simple cancellation cleanup handlers (whose only use is for  
 6657 cleaning up when a thread is canceled) is to define a general exception package that could be  
 6658 used for handling and cleaning up after hardware traps and software detected errors. This  
 6659 was too far removed from the charter of providing threads to handle asynchrony. However,  
 6660 it is an explicit goal of IEEE Std. 1003.1-200x to be compatible with existing exception  
 6661 facilities and languages having exceptions.

6662 The interaction of this facility and other procedure-based or language-level exception  
 6663 facilities is unspecified in this version of IEEE Std. 1003.1-200x. However, it is intended that it  
 6664 be possible for an implementation to define the relationship between these cancellation  
 6665 cleanup handlers and Ada, C++, or other language-level exception handling facilities.

6666 It was suggested that the cancellation cleanup handlers should also be called when the  
 6667 process exits or calls the *exec* function. This was rejected partly due to the performance  
 6668 problem caused by having to call the cancellation cleanup handlers of every thread before the  
 6669 operation could continue. The other reason was that the only state expected to be cleaned up  
 6670 by the cancellation cleanup handlers would be the intraprocess state. Any handlers that are to  
 6671 clean up the interprocess state would be registered with *atexit()*. There is the orthogonal  
 6672 problem that the *exec* functions do not honor the *atexit()* handlers, but resolving this is  
 6673 beyond the scope of IEEE Std. 1003.1-200x.

#### 6674 • Async-Cancel Safety

6675 A function is said to be *async-cancel safe* if it is written in such a way that entering the function  
 6676 with asynchronous cancelability enabled will not cause any invariants to be violated, even if  
 6677 a cancelation request is delivered at any arbitrary instruction. Functions that are *async-*  
 6678 *cancel-safe* are often written in such a way that they need to acquire no resources for their  
 6679 operation and the visible variables that they may write are strictly limited.

6680 Any routine that gets a resource as a side-effect cannot be made *async-cancel-safe* (for  
 6681 example, *malloc()*). If such a routine were called with asynchronous cancelability enabled, it  
 6682 might acquire the resource successfully, but as it was returning to the client, it could act on a  
 6683 cancelation request. In such a case, the application would have no way of knowing whether  
 6684 the resource was acquired or not.

6685 Indeed, because many interesting routines cannot be made *async-cancel-safe*, most library  
 6686 routines in general are not *async-cancel-safe*. Every library routine should specify whether or  
 6687 not it is *async-cancel safe* so that programmers know which routines can be called from code  
 6688 that is asynchronously cancelable.

### 6689 B.2.9.6 Thread Read-Write Locks

#### 6690 Background

6691 Read-write locks are often used to allow parallel access to data on multi-processors, to avoid  
 6692 context switches on uni-processors when multiple threads access the same data, and to protect  
 6693 data structures that are frequently accessed (that is, read) but rarely updated (that is, written).  
 6694 The in-core representation of a file system directory is a good example of such a data structure.  
 6695 One would like to achieve as much concurrency as possible when searching directories, but limit  
 6696 concurrent access when adding or deleting files.

6697 Although read-write locks can be implemented with mutexes and condition variables, such  
 6698 implementations are significantly less efficient than is possible. Therefore, this synchronization  
 6699 primitive is included in IEEE Std. 1003.1-200x for the purpose of allowing more efficient  
 6700 implementations in multi-processor systems.

#### 6701 Queuing of Waiting Threads

6702 The *pthread\_rwlock\_unlock()* function description states that one writer or one or more readers  
 6703 shall acquire the lock if it is no longer held by any thread as a result of the call. However, the  
 6704 function does not specify which thread(s) acquire the lock, unless the Thread Execution  
 6705 Scheduling option is supported.

6706 The standard developers considered the issue of scheduling with respect to the queuing of  
 6707 threads blocked on a read-write lock. The question turned out to be whether  
 6708 IEEE Std. 1003.1-200x should require priority scheduling of read-write locks for threads whose  
 6709 execution scheduling policy is priority-based (for example, *SCHED\_FIFO* or *SCHED\_RR*). There  
 6710 are tradeoffs between priority scheduling, the amount of concurrency achievable among readers,  
 6711 and the prevention of writer and/or reader starvation.

6712 For example, suppose one or more readers hold a read-write lock and the following threads  
 6713 request the lock in the listed order:

```
6714 pthread_rwlock_wrlock() - Low priority thread writer_a
6715 pthread_rwlock_rdlock() - High priority thread reader_a
6716 pthread_rwlock_rdlock() - High priority thread reader_b
6717 pthread_rwlock_rdlock() - High priority thread reader_c
```

6718 When the lock becomes available, should *writer\_a* block the high priority readers? Or, suppose a  
6719 read-write lock becomes available and the following are queued:

6720 pthread\_rwlock\_rdlock() - Low priority thread reader\_a  
6721 pthread\_rwlock\_rdlock() - Low priority thread reader\_b  
6722 pthread\_rwlock\_rdlock() - Low priority thread reader\_c  
6723 pthread\_rwlock\_wrlock() - Medium priority thread writer\_a  
6724 pthread\_rwlock\_rdlock() - High priority thread reader\_d

6725 If priority scheduling is applied then *reader\_d* would acquire the lock and *writer\_a* would block  
6726 the remaining readers. But should the remaining readers also acquire the lock to increase  
6727 concurrency? The solution adopted takes into account that when the Thread Execution  
6728 Scheduling option is supported, high priority threads may in fact starve low priority threads (the  
6729 application developer is responsible in this case to design the system in such a way that this  
6730 starvation is avoided). Therefore, IEEE Std. 1003.1-200x specifies that high priority readers take  
6731 precedence over lower priority writers. However, to prevent writer starvation from threads of  
6732 the same or lower priority, writers take precedence over readers of the same or lower priority.

6733 Priority inheritance mechanisms are non-trivial in the context of read-write locks. When a high  
6734 priority writer is forced to wait for multiple readers, for example, it is not clear which subset of  
6735 the readers should inherit the writer's priority. Furthermore, the internal data structures that  
6736 record the inheritance must be accessible to all readers, and this implies some sort of  
6737 serialization that could negate any gain in parallelism achieved through the use of multiple  
6738 readers in the first place. Finally, existing practice does not support the use of priority  
6739 inheritance for read-write locks. Therefore, no specification of priority inheritance or priority  
6740 ceiling is attempted. If reliable priority-scheduled synchronization is absolutely required, it can  
6741 always be obtained through the use of mutexes.

#### 6742 **Comparison to *fcntl()* Locks**

6743 The read-write locks and the *fcntl()* locks in IEEE Std. 1003.1-200x share a common goal:  
6744 increasing concurrency among readers, thus increasing throughput and decreasing delay.

6745 However, the read-write locks have two features not present in the *fcntl()* locks. First, under  
6746 priority scheduling, read-write locks are granted in priority order. Second, also under priority  
6747 scheduling, writer starvation is prevented by giving writers preference over readers of equal or  
6748 lower priority.

6749 Also, read-write locks can be used in systems lacking a file system, such as those conforming to  
6750 the minimal realtime system profile of IEEE Std. 1003.13-1998.

#### 6751 **History of Resolution Issues**

6752 Based upon some balloting objections, the draft specified the behavior of threads waiting on a  
6753 read-write lock during the execution of a signal handler, as if the thread had not called the lock  
6754 operation. However, this specified behavior would require implementations to establish  
6755 internal signal handlers even though this situation would be rare, or never happen for many  
6756 programs. This would introduce an unacceptable performance hit in comparison to the little  
6757 additional functionality gained. Therefore, the behavior of read-write locks and signals was  
6758 reverted back to its previous mutex-like specification.

6759 *B.2.9.7 Thread Interactions with Regular File Operations*

6760 There is no additional rationale for this section.

6761 **B.2.10 Sockets**

6762 There is no additional rationale for this section.

6763 *B.2.10.1 Protocol Families*

6764 There is no additional rationale for this section.

6765 *B.2.10.2 Protocols*

6766 There is no additional rationale for this section.

6767 *B.2.10.3 Addressing*

6768 There is no additional rationale for this section.

6769 *B.2.10.4 Routing*

6770 There is no additional rationale for this section.

6771 *B.2.10.5 Interfaces*

6772 There is no additional rationale for this section.

6773 *B.2.10.6 Socket Types*

6774 There is no additional rationale for this section.

6775 *B.2.10.7 Socket I/O Mode*

6776 There is no additional rationale for this section.

6777 *B.2.10.8 Socket Owner*

6778 There is no additional rationale for this section.

6779 *B.2.10.9 Socket Queue Limits*

6780 There is no additional rationale for this section.

6781 *B.2.10.10 Pending Error*

6782 There is no additional rationale for this section.

6783 *B.2.10.11 Socket Receive Queue*

6784 There is no additional rationale for this section.

6785 *B.2.10.12 Socket Out-of-Band Data State*

6786 There is no additional rationale for this section.

6787 *B.2.10.13 Connection Indication Queue*

6788           There is no additional rationale for this section.

6789 *B.2.10.14 Signals*

6790           There is no additional rationale for this section.

6791 *B.2.10.15 Asynchronous Errors*

6792           There is no additional rationale for this section.

6793 *B.2.10.16 Use of Options*

6794           There is no additional rationale for this section.

6795 *B.2.10.17 Use of Sockets for Local UNIX Connections*

6796           There is no additional rationale for this section.

6797 *B.2.10.18 Use of Sockets over Internet Protocols Based on IPv4*

6798           There is no additional rationale for this section.

6799 *B.2.10.19 Use of Sockets over Internet Protocols Based on IPv6*

6800           There is no additional rationale for this section.

6801 **B.2.11 Tracing**

6802           The organization of the tracing rationale differs from the traditional rationale in that this tracing  
6803           rationale text is written against the trace interface as a whole, rather than against the individual  
6804           components of the trace interface or the normative section in which those components are  
6805           defined. Therefore the sections below do not parallel the sections of normative text in  
6806           IEEE Std. 1003.1-200x.

6807 *B.2.11.1 Objectives*

6808           The intended uses of tracing are application-system debugging during system development, as a  
6809           “flight recorder” for maintenance of fielded systems, and as a performance measurement tool. In  
6810           all of these intended uses, the vendor-supplied computer system and its software are, for this  
6811           discussion, assumed error-free; the intent being to debug the user-written and/or third-party  
6812           application code, and their interactions. Clearly, problems with the vendor-supplied system and  
6813           its software will be uncovered from time to time, but this is a byproduct of the primary activity,  
6814           debugging user code.

6815           Another need for defining a trace interface in POSIX stems from the objective to provide an  
6816           efficient portable way to perform benchmarks. Existing practice shows that such interfaces are  
6817           commonly used in a variety of systems but with little commonality. As part of the benchmarking  
6818           needs, we must consider two aspects within the trace interface.

6819           The first, and perhaps more important one, is the qualitative aspect.

6820           The second is the quantitative aspect.

## 6821           • Qualitative Aspect

6822           To better understand this aspect, let us consider an example. Suppose that you want to  
6823           organize a number of actions to be performed during the day. Some of these actions are

6824 known at the beginning of the day. Some others, which may be more or less important, will  
 6825 be triggered by reading your mail. During the day you will make some phone calls and  
 6826 synchronously receive some more information. Finally you will receive asynchronous phone  
 6827 calls that also will trigger actions. If you, or somebody else, examines your day at work, you,  
 6828 or he, can discover that you have not efficiently organized your work. For instance, relative  
 6829 to the phone calls you made, would it be preferable to make some of these early in the  
 6830 morning? Or to delay some others until the end of the day? Relative to the phone calls you  
 6831 have received, you might find that somebody you called in the morning has called you 10  
 6832 times while you were performing some important work. To examine, afterwards, your day at  
 6833 work, you record in sequence all the trace events relative to your work. This should give you  
 6834 a chance of organizing your next day at work.

6835 This is the qualitative aspect of the trace interface. The user of a system needs to keep a trace  
 6836 of particular points the application passes through, so that he can eventually make some  
 6837 changes in the application and/or system configuration, to give the application a chance of  
 6838 running more efficiently.

6839 • Quantitative Aspect

6840 This aspect concerns primarily realtime applications, where missed deadlines can be  
 6841 undesirable. Although there are, in POSIX.1b and POSIX.1c/POSIX.1d/POSIX.1j, some  
 6842 interfaces useful for such applications (timeouts, execution time monitoring, and so on),  
 6843 there are no APIs to aid in the tuning of a realtime application's behavior (**timespec** in  
 6844 timeouts, length of message queues, duration of driver interrupt service routine, and so on).  
 6845 The tuning of an application needs a means of recording timestamped important trace events  
 6846 during execution in order to analyze offline, and eventually, to tune some realtime features  
 6847 (redesign the system with less functionalities, readjust timeouts, redesign driver interrupts,  
 6848 and so on).

6849 **Detailed Objectives**

6850 Objectives were defined to build the trace interface and are kept for historical interest. Although  
 6851 some objectives are not fully respected in this trace interface, the concept of the POSIX trace  
 6852 interface assumes the following points:

- 6853 1. It shall be possible to trace both system and user trace events concurrently.
- 6854 2. It must be possible to trace per-process trace events and also to trace system trace events  
 6855 which are unrelated to any particular process. A per-process trace event is either user-  
 6856 initiated or system-initiated.
- 6857 3. It must be possible to control tracing on a per process basis from either inside or outside  
 6858 the process.
- 6859 4. It must be possible to control tracing on a per-thread basis from inside the enclosing  
 6860 process.
- 6861 5. Trace points shall be controllable by trace event type ID from inside and outside of the  
 6862 process. Multiple trace points can have the same trace event type ID, and will be controlled  
 6863 jointly.
- 6864 6. Recording of trace events is dependent on both trace event type ID and the  
 6865 process/thread. Both must be enabled in order to record trace events. System trace events  
 6866 may or may not be handled differently.
- 6867 7. The API shall not mandate the ability to control tracing for more than one process at the  
 6868 same time.



- 6869 8. There is no objective for trace control on anything bigger than a process; for example,  
6870 group or session.
- 6871 9. Trace propagation and control:
- 6872 a. Trace propagation across fork is optional; the default is to not trace a child process.
- 6873 b. Trace control shall span *thread\_create* operations; that is, if a process is being traced,  
6874 any thread will be traced as well if this thread allows tracing. The default is to allow  
6875 tracing.
- 6876 10. Trace control shall not span *exec* or *spawn* operations.
- 6877 11. A triggering API is not required. The triggering API is the ability to command or stop  
6878 tracing based on the occurrence of specific trace event other than a `POSIX_TRACE_START`  
6879 trace event or a `POSIX_TRACE_STOP` trace event.
- 6880 12. Trace log entries shall have timestamps of implementation-defined resolution.  
6881 Implementations are exhorted to support at least microsecond resolution. When a trace log  
6882 entry is retrieved, it shall have timestamp, PC address, PID, and TID of the entity that  
6883 generated the trace event.
- 6884 13. Independently developed code should be able to use trace facilities without coordination  
6885 and without conflict.
- 6886 14. Even if the trace points in the trace calls are not unique, the trace log entries (after any  
6887 processing) shall be uniquely identified as to trace point.
- 6888 15. There shall be a standard API to read the trace stream.
- 6889 16. The format of the trace stream and the trace log is opaque and unspecified.
- 6890 17. It shall be possible to read a completed trace, if recorded on some suitable non-volatile  
6891 storage, even subsequent to a power cycle or subsequent cold boot of the system.
- 6892 18. Support of analysis of a trace log while it is being formed is implementation-defined.
- 6893 19. The API shall allow the application to write trace stream identification information into the  
6894 trace stream and to be able to retrieve it, without it being overwritten by trace entries, even  
6895 if the trace stream is full.
- 6896 20. It must be possible to specify the destination of trace data produced by trace events.
- 6897 21. It must be possible to have different trace streams, and for the tracing enabled by one trace  
6898 stream to be completely independent of the tracing of another trace stream.
- 6899 22. It must be possible to trace events from threads in different CPUs.
- 6900 23. The API shall support one or more trace streams per-system, and one or more trace  
6901 streams per-process, up to an implementation-defined set of per-system and per-process  
6902 maximums.
- 6903 24. It shall be possible to determine the order in which the trace events happened, without  
6904 necessarily depending on the clock, up to an implementation-defined time resolution.
- 6905 25. For performance reasons, the trace event point call(s) shall be implementable as a macro  
6906 (see the ISO POSIX-1: 1996 standard, Subclause 1.3.4, Statement 2).
- 6907 26. IEEE Std. 1003.1-200x must not define the trace points which a conforming system must  
6908 implement, except for trace points used in the control of tracing.
- 6909 27. The APIs shall be thread-safe, and trace points should be lock-free (that is shall not require  
6910 a lock to gain exclusive access to some resource).

- 6911 28. The user-provided information associated with a trace event is variable-sized, up to some  
6912 maximum size.
- 6913 29. Bounds on record and trace stream sizes:
- 6914 a. The API must permit the application to declare the upper bounds on the length of an  
6915 application data record. The system shall return the limit it used. The limit used may  
6916 be smaller than requested.
- 6917 b. The API must permit the application to declare the upper bounds on the size of trace  
6918 streams. The system shall return the limit it used. The limit used may be different,  
6919 either larger or smaller, than requested.
- 6920 30. The API must be able to pass any fundamental data type, and a structured data type  
6921 composed only of fundamental types. The API must be able to pass data by reference,  
6922 given only as an address and a length. Fundamental types are the POSIX.1 types (see the  
6923 ISO POSIX-1:1996 standard, Subclause 2.5, Table 2-1) plus those defined in the ISO C  
6924 standard.
- 6925 31. The API shall apply the POSIX notions of ownership and permission to recorded trace  
6926 data, corresponding to the sources of that data.

### 6927 **Comments on Objectives**

- 6928 **Note:** In the following comments, numbers in square brackets refer to the above objectives.
- 6929 It is necessary to be able to obtain a trace stream for a complete activity. This means we need to  
6930 be able to trace both application and system trace events. A per-process trace event is either  
6931 user-initiated, like the *write()* POSIX call, or system-initiated, like a timer expiration. We also  
6932 need to be able to trace an entire process's activity even when it has threads in multiple CPUs. To  
6933 avoid excess trace activity, it is necessary to be able to control tracing on a trace event type basis.  
6934 [Objectives 1,2,5,22]
- 6935 We need to be able to control tracing on a per-process basis, both from inside and outside the  
6936 process; that is, a process can start a trace activity on itself or any other process. We also see the  
6937 need to allow the definition of a maximum number trace streams per system.  
6938 [Objectives 3,23]
- 6939 From within a process, it is necessary to be able to control tracing on a per-thread basis. This  
6940 provides an additional filtering capability to keep the amount of traced data to a minimum. It  
6941 also allows for less ambiguity as to the origin of trace events. It is recognized that thread-level  
6942 control is only valid from within the process itself. It is also desirable to know the maximum  
6943 number of trace streams per process that can be started. We do not want the API to require  
6944 thread synchronization or to mandate priority inversions that would cause the thread to block.  
6945 However, the API must be thread-safe.  
6946 [Objectives 4,23,24,27]
- 6947 We see no objective to control tracing on anything larger than a process; for example, a group or  
6948 session. Also, the ability to start or stop a trace activity on multiple processes atomically may be  
6949 very difficult or cumbersome in some implementations.  
6950 [Objectives 6,8]
- 6951 It is also necessary to be able to control tracing by trace event type identifier, sometimes called a  
6952 trace hook ID. However, there is no mandated set of system trace events, since such trace points  
6953 are very system-dependent. The API must not require from the operating system facilities that  
6954 are not standard (POSIX).  
6955 [Objectives 6,26]

6956 Trace control must span *fork()* and *pthread\_create()*. If not, there will be no way to ensure that a  
6957 program's activity is entirely traced. The newly forked child would not be able to turn on its  
6958 tracing until after it obtained control after the fork, and trace control externally would be even  
6959 more problematic.  
6960 [Objective 9]

6961 Since *exec()* and *spawn()* represent a complete change in the execution of a task (a new  
6962 program), trace control need not persist over an *exec()* or *spawn()*.  
6963 [Objective 10]

6964 Where trace activities are started on multiple processes, these trace activities should not interfere  
6965 with each other.  
6966 [Objective 21]

6967 There is no need for a triggering objective, primarily for performance reasons; see also Section  
6968 B.2.11.8 (on page 3498), rationale on triggering.  
6969 [Objective 11]

6970 It must be possible to determine the origin of each traced event. We need the process and thread  
6971 identifiers for each trace event. We also saw the need for a user-specifiable origin, but felt this  
6972 would create too much overhead.  
6973 [Objectives 12,14]

6974 We must allow for trace points to come embedded in software components from several  
6975 different sources and vendors without requiring coordination.  
6976 [Objective 13]

6977 We need to be able to uniquely identify trace points that may have the same trace stream  
6978 identifier. We only need to be able to do this when a trace report is produced.  
6979 [Objectives 12,14]

6980 Tracing is a very performance-sensitive activity, and will therefore likely be implemented at a  
6981 low level within the system. Hence the interface shall not mandate any particular buffering or  
6982 storage method. Therefore, we will need a standard API to read a trace stream. Also the interface  
6983 shall not mandate the format of the trace data, and the interface shall not assume a trace storage  
6984 method. Due to the possibility of a monolithic kernel and the possible presence of multiple  
6985 processes capable of running trace activities, the two kinds of trace events may be stored in two  
6986 separate streams for performance reasons. A mandatory dump mechanism, common in some  
6987 existing practice, has been avoided to allow the implementation of this set of functions on small  
6988 realtime profiles for which the concept of a file system is not defined. The trace API calls should  
6989 be implemented as macros.  
6990 [Objectives 15,16,25,30]

6991 Since a trace facility is a valuable service tool, the output (or log) of a completed trace stream  
6992 that is written to permanent storage must be readable on other systems of the type that  
6993 produced the trace log. Note that there is no objective to be able to interpret a trace log that was  
6994 not successfully completed.  
6995 [Objectives 17,18,19]

6996 For trace streams written to permanent storage, a way to specify the destination of the trace  
6997 stream is needed.  
6998 [Objective 20]

6999 We need to be able to depend on the ordering of trace events up to some system-defined time  
7000 interval. For example, we need to know the time period which, if trace events are closer together,  
7001 their ordering is indeterminate. Events that occur within an interval smaller than this resolution  
7002 may or may not be read back in the correct order.

- 7003 [Objective 24]
- 7004 The application should be able to know how much data can be traced. When trace event types
- 7005 can be filtered, the application should be able to specify the approximate maximum amount of
- 7006 data that will be traced in a trace event so resources can be more efficiently allocated.
- 7007 [Objectives 28,29]
- 7008 Users should not be able to trace data to which they would not normally have access to. System
- 7009 trace events corresponding to a process/thread should be associated with the ownership of that
- 7010 process/thread.
- 7011 [Objective 31]

7012 *B.2.11.2 Trace Model*

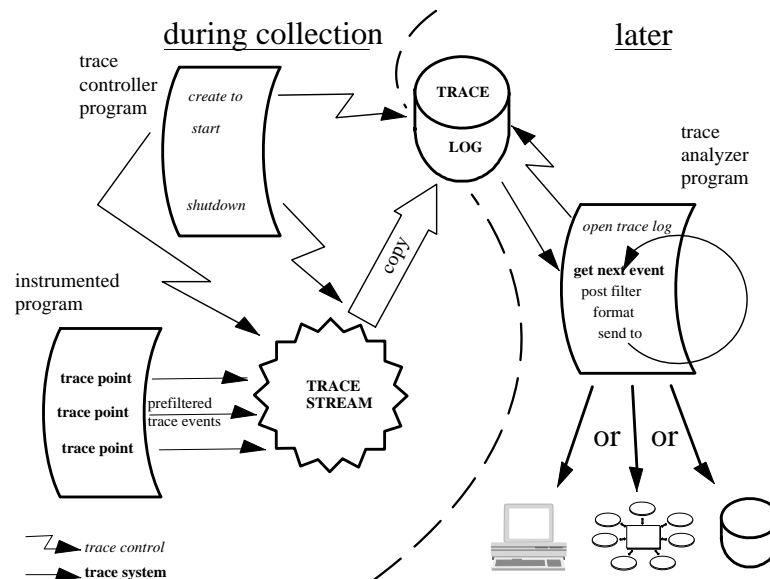
7013 **Introduction**

7014 The model is based on two base entities: the “Trace Stream” and the “Trace Log” , and a  
 7015 recorded unit called the “Trace Event”. The possibility of using Trace Streams and Trace Logs  
 7016 separately gives us two use dimensions and solves both the performance issue and the full-  
 7017 information system issue. In the case of a trace stream without log, specific information,  
 7018 although reduced in quantity, is required to be registered, in a possibly small realtime system,  
 7019 with as little overhead as possible. The Trace Log option has been added for small realtime  
 7020 systems. In the case of a trace stream with log, considerable complex application-specific  
 7021 information needs to be collected.

7022 **Trace Model Description**

7023 The trace model can be examined for three different subfunctions: Application Instrumentation,  
 7024 Trace Operation Control, and Trace Analysis.

7025



7026 **Figure B-2** Trace System Overview: for Offline Analysis

7027 Each of these subfunctions requires specific characteristics of the trace mechanism API.

## 7028 • Application Instrumentation

7029 When instrumenting an application, the programmer has no concern about the future  
 7030 utilization of the trace events in trace stream or trace log, the full policy of trace stream, or  
 7031 the eventual pre-filtering of trace events. But he is concerned about the correct determination  
 7032 of specific trace event type identifier, regardless of how many independent libraries are used  
 7033 in the same user application; see Figure B-2 and Figure B-3 (on page 3482).

7034 This trace API shall provide the necessary operations to accomplish this subfunction. This is  
 7035 done by providing functions to associate a programmer-defined name with an  
 7036 implementation-defined trace event type identifier; see the *posix\_trace\_eventid\_open()*  
 7037 function), and to send this trace event into a potential trace stream (see the  
 7038 *posix\_trace\_event()* function).

## 7039 • Trace Operation Control

7040 When controlling the recording of trace events in a trace stream, the programmer is  
 7041 concerned with the correct initialization of the trace mechanism (that is, the sizing of the  
 7042 trace stream), the correct retention of trace events in a permanent storage, the correct  
 7043 dynamic recording of trace events, and so on.

7044 This trace API shall provide the necessary material to permit this efficiently. This is done by  
 7045 providing functions to initialize a new trace stream, and optionally a trace log:

- 7046 — Trace Stream Attributes Object Initialization (see *posix\_trace\_attr\_init()*)
- 7047 — Functions to Retrieve or Set Information About a Trace Stream (see  
 7048 *posix\_trace\_attr\_getgenversion()*)
- 7049 — Functions to Retrieve or Set the Behavior of a Trace Stream (see  
 7050 *posix\_trace\_attr\_getinherited()*)
- 7051 — Functions to Retrieve or Set Trace Stream Size Attributes (see  
 7052 *posix\_trace\_attr\_getmaxusevents()*)
- 7053 — Trace Stream Initialization, Flush, and Shutdown from a Process (see *posix\_trace\_create()*)
- 7054 — Clear Trace Stream and Trace Log (see *posix\_trace\_clear()*)

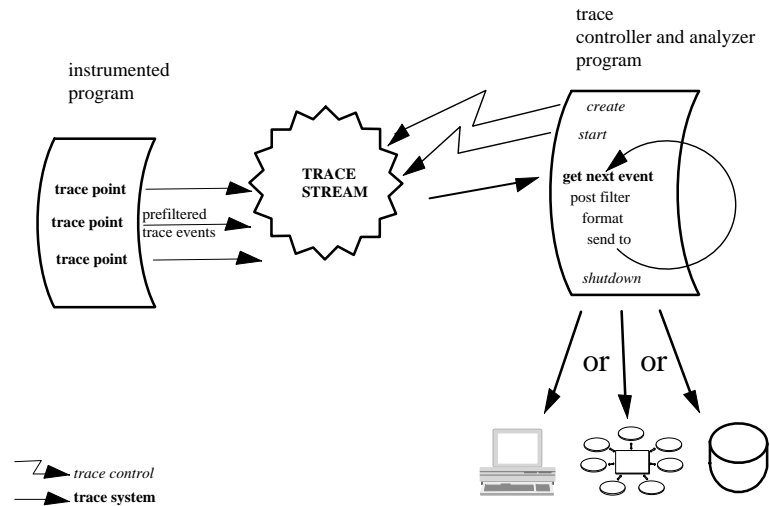
7055 To select the trace event types that are to be traced:

- 7056 — Manipulate Trace Event Type Identifier (see *posix\_trace\_trid\_eventid\_open()*)
- 7057 — Iterate over a Mapping of Trace Event Type (see *posix\_trace\_eventtypelist\_getnext\_id()*)
- 7058 — Manipulate Trace Event Type Sets (see *posix\_trace\_eventset\_empty()*)
- 7059 — Set Filter of an Initialized Trace Stream (see *posix\_trace\_set\_filter()*)

7060 To control the execution of an active trace stream:

- 7061 — Trace Start and Stop (see *posix\_trace\_start()*)
- 7062 — Functions to Retrieve the Trace Attributes or Trace Statuses (see *posix\_trace\_get\_attr()*)

7063



7064

**Figure B-3** Trace System Overview: for Online Analysis

7065

- Trace Analysis

7066

7067

7068

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.

7069

7070

7071

7072

7073

7074

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.

7075

7076

7077

7078

7079

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.

7080

This trace API shall provide functions to:

7081

- Extract trace stream identification attributes (see *posix\_trace\_attr\_getgenversion()*)

7082

- Extract trace stream behavior attributes (see *posix\_trace\_attr\_getinherited()*)

7083

7084

- Extract trace event, stream, and log size attributes (see *posix\_trace\_attr\_getmaxusereventsize()*)

7085

- Look up trace event type names (see *posix\_trace\_eventid\_get\_name()*)

- 7086 — Iterate over trace event type identifiers (see *posix\_trace\_eventtypelist\_getnext\_id()*)
- 7087 — Open, rewind, and close a trace log (see *posix\_trace\_open()*)
- 7088 — Read trace stream attributes and status (see *posix\_trace\_get\_attr()*)
- 7089 — Read trace events (see *posix\_trace\_getnext\_event()*)

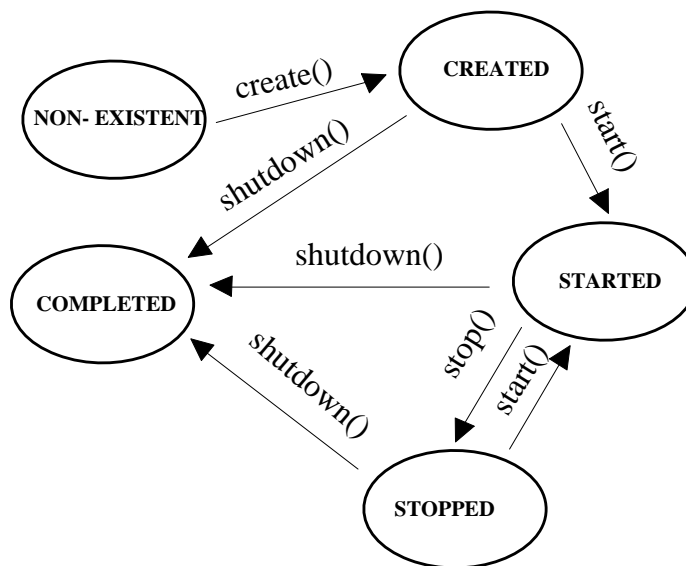
7090 Due to the following two reasons:

- 7091 1. The requirement that the trace system must not add unacceptable overhead to the traced process and so that the trace event point execution must be fast
- 7092
- 7093 2. The traced application does not care about tracing errors

7094 the trace system cannot return any internal error to the application. Internal error conditions can range from unrecoverable errors that will force the active trace stream to abort, to small errors that can affect the quality of tracing without aborting the trace stream. The group decided to define a system trace event to report to the analysis process such internal errors. It is not the intention of IEEE Std. 1003.1-200x to require an implementation to report an internal error that corrupts or terminates tracing operation. The implementor is free to decide which internal documented errors, if any, the trace system is able to report.

7101 **States of a Trace Stream**

7102



7103 **Figure B-4** Trace System Overview: States of a Trace Stream

7104 Figure B-4 shows the different states an active trace stream passes through. After the  
 7105 *posix\_trace\_create()* function call, a trace stream becomes CREATED and a trace stream is  
 7106 associated for the future collection of trace events. The status of the trace stream is  
 7107 POSIX\_TRACE\_SUSPENDED. The state becomes STARTED after a call to the *posix\_trace\_start()*  
 7108 function, and the status becomes POSIX\_TRACE\_RUNNING. In this state, all trace events that  
 7109 are not filtered out shall be stored into the trace stream. After a call to *posix\_trace\_stop()*, the

7110 trace stream becomes STOPPED (and the status POSIX\_TRACE\_SUSPENDED). In this state, no  
 7111 new trace events will be recorded in the trace stream, but previously recorded trace events may  
 7112 continue to be read.

7113 After a call to *posix\_trace\_shutdown()*, the trace stream is in the state COMPLETED. The trace  
 7114 stream no longer exists but, if the Trace Log option is supported, all the information contained in  
 7115 it has been logged. If a log object has not been associated with the trace stream at the creation, it  
 7116 is the responsibility of the trace controller process to not shut the trace stream down while trace  
 7117 events remain to be read in the stream.

### 7118 **Tracing All Processes**

7119 Some implementations have a tracing subsystem with the ability to trace all processes. This is  
 7120 useful to debug some types of device drivers such as those for ATM or X25 adapters. These types  
 7121 of adapters are used by several independent processes, that are not issued from the same  
 7122 process.

7123 The POSIX trace interface does not define any constant or option to create a trace stream tracing  
 7124 all processes. But the POSIX trace interface does not prevent this type of implementation and the  
 7125 implementor is free to add this capability. Nevertheless, the POSIX trace interface allows to trace  
 7126 all the system trace events and all the processes issued from the same process.

7127 If such a tracing system capability has to be implemented, when a trace stream is created, it is  
 7128 recommended that a constant named POSIX\_TRACE\_ALLPROC be used instead of the process  
 7129 identifier in the argument of the function *posix\_trace\_create()* or *posix\_trace\_create\_withlog()*. A  
 7130 possible value for POSIX\_TRACE\_ALLPROC may be -1 instead of a real process identifier.

7131 The implementor has to be aware that there is some impact on the tracing behavior as defined in  
 7132 the POSIX trace interface. For example:

- 7133 • If the default value for the inheritance attribute is to set to  
 7134 POSIX\_TRACE\_CLOSE\_FOR\_CHILD, the implementation has to stop tracing for the child  
 7135 process.
- 7136 • The trace controller which is creating this type of trace stream must have the appropriate  
 7137 privilege to trace all the processes.

### 7138 **Trace Storage**

7139 The model is based on two types of trace events: system trace events and user-defined trace  
 7140 events. The internal representation of trace events is implementation-defined, and so the  
 7141 implementor is free to choose the more suitable, practical, and efficient way to design the  
 7142 internal management of trace events. For the timestamping operation, the model does not  
 7143 impose the CLOCK\_REALTIME or any other clock. The buffering allocation and operation  
 7144 follow the same principle. The implementor is free to use one or more buffers to record trace  
 7145 events; the interface assumes only a logical trace stream of sequentially recorded trace events.  
 7146 Regarding flushing of trace events, the interface allows the definition of a trace log object which  
 7147 typically can be a file. But the group was also aware of defining functions to permit the use of  
 7148 this interface in small realtime systems, which may not have general file system capabilities. For  
 7149 instance, the three functions *posix\_trace\_getnext\_event()* (blocking),  
 7150 *posix\_trace\_timedgetnext\_event()* (blocking with timeout), and *posix\_trace\_trygetnext\_event()*  
 7151 (non-blocking) are proposed to read the recorded trace events.

7152 The policy to be used when the trace stream becomes full also relies on common practice:

- 7153 • For an active trace stream, the POSIX\_TRACE\_LOOP trace stream policy permits automatic  
 7154 overrun (overwrite of oldest trace events) while waiting for some user-defined condition to



7155 cause tracing to stop. By contrast, the POSIX\_TRACE\_UNTIL\_FULL trace stream policy  
 7156 requires the system to stop tracing when the trace stream is full. However, if the trace stream  
 7157 that is full is at least partially emptied by a call to the *posix\_trace\_flush()* function or by calls  
 7158 to *posix\_trace\_getnext\_event()* function, the trace system will automatically resume tracing.

7159 If the Trace Log option is supported the operation of the POSIX\_TRACE\_FLUSH policy is an  
 7160 extension of the POSIX\_TRACE\_UNTIL\_FULL policy. The automatic free operation (by  
 7161 flushing to the associated trace log) is added.

7162 • If a log is associated with the trace stream and this log is a regular file, these policies also  
 7163 apply for the log. One more policy, POSIX\_TRACE\_APPEND, is defined to allow indefinite  
 7164 extension of the log. Since the log destination can be any device or pseudo-device, the  
 7165 implementation may not be able to manipulate the destination as required by  
 7166 IEEE Std. 1003.1-200x. For this reason, the behavior of the log full policy may be unspecified  
 7167 depending of the trace log type.

7168 The current trace interface does not define a service to preallocate space for a trace log file,  
 7169 because this space can be preallocated by means of a call to the *posix\_fallocate()* function. This  
 7170 function could be called after the file has been opened, but before the trace stream is created.  
 7171 The *posix\_fallocate()* function ensures that any required storage for regular file data is  
 7172 allocated on the file system storage media. If *posix\_fallocate()* returns successfully,  
 7173 subsequent writes to the specified file data shall not fail due to the lack of free space on the  
 7174 file system storage media. Besides trace events, a trace stream also includes trace attributes  
 7175 and the mapping from trace event names to trace event type identifiers. The implementor is  
 7176 free to choose how to store the trace attributes and the trace event type map, but must ensure  
 7177 that this information is not lost when a trace stream overrun occurs.

### 7178 B.2.11.3 Trace Programming Examples

7179 Several programming examples are presented to show the code of the different possible  
 7180 subfunctions using a trace subsystem. All these programs need to include the `<trace.h>` header.  
 7181 In the examples shown, error checking is omitted for more simplicity.

### 7182 Trace Operation Control

7183 These examples show the creation of a trace stream for another process; one which is already  
 7184 trace instrumented. All the default trace stream attributes are used to simplify programming in  
 7185 the first example. The second example shows more possibilities.

### 7186 First Example

```
7187 /* Caution. Error checks omitted */
7188 {
7189     trace_attr_t attr;
7190     pid_t pid = traced_process_pid;
7191     int fd;
7192     trace_id_t trid;
7193     - - - - -
7194     /* Initialize trace stream attributes */
7195     posix_trace_attr_init(&attr);
7196     /* Open a trace log */
7197     fd=open("/tmp/mytracelog",...);
7198     /*
7199      * Create a new trace associated with a log
7200      * and with default attributes
```

```

7201         */
7202         posix_trace_create_withlog(pid, &attr, fd, &trid);
7203         /* Trace attribute structure can now be destroyed */
7204         posix_trace_attr_destroy(&attr);
7205         /* Start of trace event recording */
7206         posix_trace_start(trid);
7207         - - - - -
7208         - - - - -
7209         /* Duration of tracing */
7210         - - - - -
7211         - - - - -
7212         /* Stop and shutdown of trace activity */
7213         posix_trace_shutdown(trid);
7214         - - - - -
7215     }

```

## 7216 **Second Example**

7217 Between the initialization of the trace stream attributes and the creation of the trace stream,  
7218 these trace stream attributes may be modified; see **Trace Stream Attribute Manipulation** (on  
7219 page 3490) for specific programming example. Between the creation and the start of the trace  
7220 stream, the event filter may be set; after the trace stream is started, the event filter may be  
7221 changed. The setting of an event set and the change of a filter is shown in **Create a Trace Event**  
7222 **Type Set and Change the Trace Event Type Filter** (on page 3490).

```

7223     /* Caution. Error checks omitted */
7224     {
7225         trace_attr_t attr;
7226         pid_t pid = traced_process_pid;
7227         int fd;
7228         trace_id_t trid;
7229         - - - - -
7230         /* Initialize trace stream attributes */
7231         posix_trace_attr_init(&attr);
7232         /* Attr default may be changed at this place; see example */
7233         - - - - -
7234         /* Create and open a trace log with R/W user access */
7235         fd=open("/tmp/mytracelog",O_WRONLY|O_CREAT,S_IRUSR|S_IWUSR);
7236         /* Create a new trace associated with a log */
7237         posix_trace_create_withlog(pid, &attr, fd, &trid);
7238         /*
7239          * If the Trace Filter option is supported
7240          * trace event type filter default may be changed at this place;
7241          * see example about changing the trace event type filter
7242          */
7243         posix_trace_start(trid);
7244         - - - - -
7245
7246         /*
7247          * If you have an uninteresting part of the application
7248          * you can stop temporarily.
7249          */

```

```

7249     * posix_trace_stop(trid);
7250     * - - - - -
7251     * - - - - -
7252     * posix_trace_start(trid);
7253     */
7254     - - - - -
7255     /*
7256     * If the Trace Filter option is supported
7257     * the current trace event type filter can be changed
7258     * at any time (see example about how to set
7259     * a trace event type filter
7260     */
7261     - - - - -
7262     /* Stop the recording of trace events */
7263     posix_trace_stop(trid);
7264     /* Shutdown the trace stream */
7265     posix_trace_shutdown(trid);
7266     /*
7267     * Destroy trace stream attributes; attr structure may have
7268     * been used during tracing to fetch the attributes
7269     */
7270     posix_trace_attr_destroy(&attr);
7271     - - - - -
7272 }

```

### 7273 Application Instrumentation

7274 This example shows an instrumented application. The code is included in a block of instructions,  
7275 perhaps a function from a library. Possibly in an initialization part of the instrumented  
7276 application, two user trace event names are mapped to two trace event type identifiers  
7277 (function *posix\_trace\_eventid\_open()*). Then two trace points are programmed.

```

7278 /* Caution. Error checks omitted */
7279 {
7280     trace_eventid_t eventid1, eventid2;
7281     - - - - -
7282     /* Initialization of two trace event type ids */
7283     posix_trace_eventid_open("my_first_event",&eventid1);
7284     posix_trace_eventid_open("my_second_event",&eventid2);
7285     - - - - -
7286     - - - - -
7287     - - - - -
7288     /* Trace point */
7289     posix_trace_event(eventid1,NULL,0);
7290     - - - - -
7291     /* Trace point */
7292     posix_trace_event(eventid2,NULL,0);
7293     - - - - -
7294 }

```

7295 **Trace Analyzer**

7296 This example shows the manipulation of a trace log resulting from the dumping of a completed  
 7297 trace stream. All the default attributes are used to simplify programming, and data associated  
 7298 with a trace event are not shown in the first example. The second example shows more  
 7299 possibilities.

7300 **First Example**

```

7301 /* Caution. Error checks omitted */
7302 {
7303     int fd;
7304     trace_id_t trid;
7305     posix_trace_event_info trace_event;
7306     char trace_event_name[TRACE_EVENT_NAME_MAX];
7307     int return_value;
7308     size_t returndatasize;
7309     int lost_event_number;
7310
7311     - - - - -
7312
7313     /* Open an existing trace log */
7314     fd=open("/tmp/tracelog", O_RDONLY);
7315     /* Open a trace stream on the open log */
7316     posix_trace_open(fd, &trid);
7317     /* Read a trace event */
7318     posix_trace_getnext_event(trid, &trace_event,
7319                             NULL, 0, &returndatasize,&return_value);
7320
7321     /* Read and print all trace event names out in a loop */
7322     while (return_value == NULL)
7323     {
7324         /*
7325          * Get the name of the trace event associated
7326          * with trid trace ID
7327          */
7328         posix_trace_eventid_get_name(trid, trace_event.event_id,
7329                                     trace_event_name);
7330         /* Print the trace event name out */
7331         printf("%s\n",trace_event_name);
7332         /* Read a trace event */
7333         posix_trace_getnext_event(trid, &trace_event,
7334                                 NULL, 0, &returndatasize,&return_value);
7335     }
7336
7337     /* Close the trace stream */
7338     posix_trace_close(trid);
7339     /* Close the trace log */
7340     close(fd);
7341 }

```

7338 **Second Example**

7339 The complete example includes the two other examples in **Retrieve Information from a Trace**  
 7340 **Log** (on page 3491) and in **Retrieve the List of Trace Event Types Used in a Trace Log** (on page  
 7341 3492). For example, the *maxdatasize* variable is set in **Retrieve the List of Trace Event Types**  
 7342 **Used in a Trace Log** (on page 3492).

```

7343 /* Caution. Error checks omitted */
7344 {
7345     int fd;
7346     trace_id_t trid;
7347     posix_trace_event_info trace_event;
7348     char trace_event_name[TRACE_EVENT_NAME_MAX];
7349     char * data;
7350     size_t maxdatasize=1024, returndatasize;
7351     int return_value;
7352     - - - - -
7353     /* Open an existing trace log */
7354     fd=open("/tmp/tracelog", O_RDONLY);
7355     /* Open a trace stream on the open log */
7356     posix_trace_open( fd, &trid);
7357     /*
7358      * Retrieve information about the trace stream which
7359      * was dumped in this trace log (see example)
7360      */
7361     - - - - -
7362     /* Allocate a buffer for trace event data */
7363     data=(char *)malloc(maxdatasize);
7364     /*
7365      * Retrieve the list of trace event used in this
7366      * trace log (see example)
7367      */
7368     - - - - -
7369     /* Read and print all trace event names and data out in a loop */
7370     while (1)
7371     {
7372     posix_trace_getnext_event(trid, &trace_event,
7373         data, maxdatasize, &returndatasize,&return_value);
7374         if (return_value != NULL) break;
7375         /*
7376          * Get the name of the trace event type associated
7377          * with trid trace ID
7378          */
7379         posix_trace_eventid_get_name(trid, trace_event.event_id,
7380             trace_event_name);
7381         {
7382         int i;
7383
7384         /* Print the trace event name out */
7385         printf("%s: ", trace_event_name);
7386         /* Print the trace event data out */
7387         for (i=0; i<returndatasize, i++) printf("%02.2X",

```

```

7387         (unsigned char)data[i]);
7388     printf("\n");
7389     }
7390 }
7391     /* Close the trace stream */
7392     posix_trace_close(trid);
7393     /* The buffer data is deallocated */
7394     free(data);
7395     /* Now the file can be closed */
7396     close(fd);
7397 }

```

### 7398 **Several Programming Manipulations**

7399 The following examples show some typical sets of operations needed in some contexts.

#### 7400 **Trace Stream Attribute Manipulation**

7401 This example shows the manipulation of a trace stream attribute object in order to change the  
7402 default value provided by a previous *posix\_trace\_attr\_init()* call.

```

7403 /* Caution. Error checks omitted */
7404 {
7405     trace_attr_t attr;
7406     size_t logsize=100000;
7407     - - - - -
7408     /* Initialize trace stream attributes */
7409     posix_trace_attr_init(&attr);
7410     /* Set the trace name in the attributes structure */
7411     posix_trace_attr_setname(&attr, "my_trace");
7412     /* Set the trace full policy */
7413     posix_trace_attr_setstreamfullpolicy(&attr, POSIX_TRACE_LOOP);
7414     /* Set the trace log size */
7415     posix_trace_attr_setlogsize(&attr, logsize);
7416     - - - - -
7417 }

```

#### 7418 **Create a Trace Event Type Set and Change the Trace Event Type Filter**

7419 This example is valid only if the Trace Event Filter option is supported. This example shows the  
7420 manipulation of a trace event type set in order to change the trace event type filter for an existing  
7421 active trace stream, which may be just-created, running, or suspended. Some sets of trace event  
7422 types are well-known, such as the set of trace event types not associated with a process, some  
7423 trace event types are just-built trace event types for this trace stream; one trace event type is the  
7424 predefined trace event error type which is deleted from the trace event type set.

```

7425 /* Caution. Error checks omitted */
7426 {
7427     trace_id_t trid = existing_trace;
7428     trace_event_set_t set;
7429     trace_event_id_t trace_event1, trace_event2;
7430     - - - - -
7431     /* Initialize to an empty set of trace event types */

```

```

7432     posix_trace_eventset_emptyset(&set);
7433     /*
7434     * Fill the set with all system trace events
7435     * not associated with a process
7436     */
7437     posix_trace_eventset_fill(&set, POSIX_TRACE_WOPID_EVENTS);
7438     /*
7439     * Get the trace event type identifier of the known trace event name
7440     * my_first_event for the trid trace stream
7441     */
7442     posix_trace_trid_eventid_open(trid, "my_first_event", &trace_event1);
7443     /* Add the set with this trace event type identifier */
7444     posix_trace_eventset_add_event(trace_event1, &set);
7445     /*
7446     * Get the trace event type identifier of the known trace event name
7447     * my_second_event for the trid trace stream
7448     */
7449
7450     posix_trace_trid_eventid_open(trid, "my_second_event", &trace_event2);
7451     /* Add the set with this trace event type identifier */
7452     posix_trace_eventset_add_event(trace_event2, &set);
7453     - - - - -
7454     /* Delete the system trace event POSIX_TRACE_ERROR from the set */
7455     posix_trace_eventset_del_event(POSIX_TRACE_ERROR, &set);
7456     - - - - -
7457     /* Modify the trace stream filter making it equal to the new set */
7458     posix_trace_set_filter(trid, &set, POSIX_TRACE_SET_EVENTSET);
7459     - - - - -
7460     /*
7461     * Now trace_event1, trace_event2, and all system trace event types
7462     * not associated with a process, except for the POSIX_TRACE_ERROR
7463     * system trace event type, are filtered out of (not recorded in) the
7464     * existing trace stream.
7465     */
7466     }

```

#### 7466 **Retrieve Information from a Trace Log**

7467 This example shows how to extract information from a trace log, the dump of a trace stream.  
7468 This code:

```

7469     • Asks if the trace stream has lost trace events
7470     • Extracts the information about the version of the trace subsystem which generated this trace
7471     log
7472     • Retrieves the maximum size of trace event data; this may be used to dynamically allocate an
7473     array for extracting trace event data from the trace log without overflow
7474     /* Caution. Error checks omitted */
7475     {
7476         struct posix_trace_status_info statusinfo;
7477         trace_attr_t attr;
7478         trace_id_t trid = existing_trace;

```

```

7479     size_t maxdatasize;
7480     char genversion[TRACE_NAME_MAX];
7481     - - - - -
7482     /* Get the trace stream status */
7483     posix_trace_get_status(trid, &statusinfo);
7484     /* Detect an overrun condition */
7485     if (statusinfo.posix_stream_overrun_status == POSIX_TRACE_OVERRUN)
7486         printf("trace events have been lost\n");
7487
7488     /* Get attributes from the trid trace stream */
7489     posix_trace_get_attr(trid, &attr);
7490     /* Get the trace generation version from the attributes */
7491     posix_trace_attr_getgenversion(&attr, genversion);
7492     /* Print the trace generation version out */
7493     printf("Information about Trace Generator:%s\n",genversion);
7494
7495     /* Get the trace event max data size from the attributes */
7496     posix_trace_attr_getmaxdatasize(&attr, &maxdatasize);
7497     /* Print the trace event max data size out */
7498     printf("Maximum size of associated data:%d\n",maxdatasize);
7499     /* Destroy the trace stream attributes */
7500     posix_trace_attr_destroy(&attr);
7501 }

```

### Retrieve the List of Trace Event Types Used in a Trace Log

This example shows the retrieval of a trace stream's trace event type list. This operation may be very useful if you are interested only in tracking the type of trace events in a trace log.

```

7503     /* Caution. Error checks omitted */
7504     {
7505         trace_id_t trid = existing_trace;
7506         trace_event_id_t event_id;
7507         char event_name[TRACE_EVENT_NAME_MAX];
7508         int return_value;
7509         - - - - -
7510
7511         /*
7512          * In a loop print all existing trace event names out
7513          * for the trid trace stream
7514          */
7515         while (1)
7516         {
7517             posix_trace_eventtypelist_getnext_id(trid, &event_id
7518                 &return_value);
7519             if (return_value != NULL) break;
7520             /*
7521              * Get the name of the trace event associated
7522              * with trid trace ID
7523              */
7524             posix_trace_eventid_get_name(trid, event_id, event_name);
7525             /* Print the name out */
7526             printf("%s\n", event_name);
7527         }

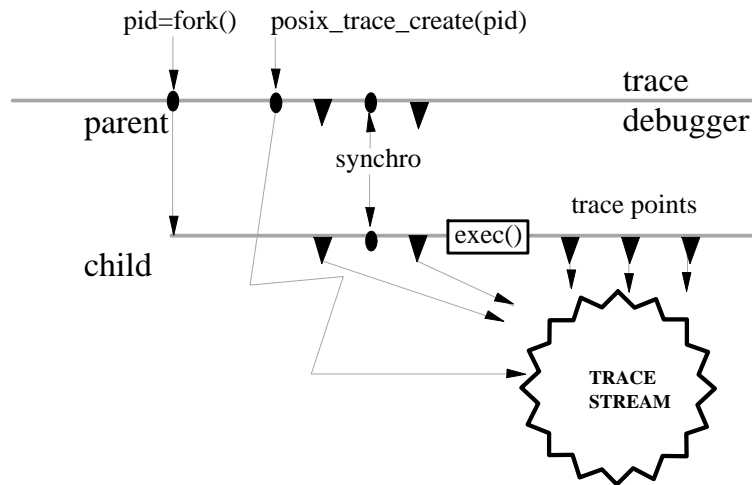
```



7527 }

7528 *B.2.11.4 Rationale on Trace for Debugging*

7529



7530

**Figure B-5** Trace Another Process

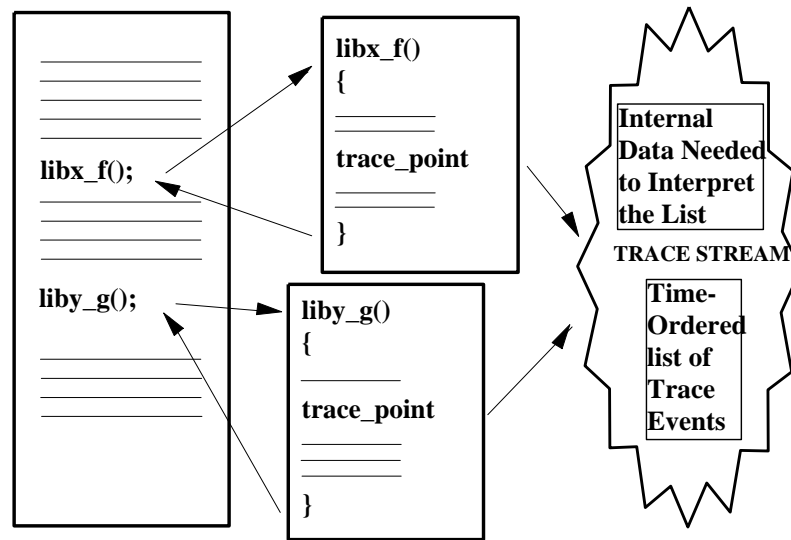
7531 Among the different possibilities offered by the trace interface defined in IEEE Std. 1003.1-200x,  
 7532 the debugging of an application is the most interesting one. Typical operations in the controlling  
 7533 debugger process are to filter trace event types, to get trace events from the trace stream, to stop  
 7534 the trace stream when the debugged process is executing uninteresting code, to start the trace  
 7535 stream when some interesting point is reached, and so on. The interface defined in  
 7536 IEEE Std. 1003.1-200x should define all the necessary base functions to allow this dynamic debug  
 7537 handling.

7538 Figure B-5 shows an example in which the trace stream is created after the call to the *fork()*  
 7539 function. If the user does not want to lose trace events some synchronization mechanism  
 7540 (represented in the figure) may be needed before calling the *exec()* function, to give the parent  
 7541 a chance to create the trace stream before the child begins the execution of its trace points.

7542 *B.2.11.5 Rationale on Trace Event Type Name Space*

7543 At first, the working group was in favor of the representation of a trace event type by an integer  
 7544 (*event\_name*). It seems that existing practice shows the weakness of such a representation. The  
 7545 collision of trace event types is the main problem that cannot be simply resolved using this sort  
 7546 of representation. Suppose, for example, that a third party designs an instrumented library. The  
 7547 user does not have the source of this library and wants to trace his application which uses in  
 7548 some part the third-party library. There is no means for him to know what are the trace event  
 7549 types used in the instrumented library so he has some chance of duplicating some of them and  
 7550 thus to obtain a contaminated tracing of his application.

7551



7552

**Figure B-6** Trace Name Space Overview: With Third-Party Library

7553

7554

7555

7556

7557

7558

7559

7560

7561

7562

7563

We have requirements to allow program images containing pieces from various vendors to be traced without also requiring those or any other vendors to coordinate their uses of the trace facility, and especially the naming of their various trace event types and trace point IDs. The chosen solution is to provide a very large name space, large enough so that the individual vendors can give their trace types and tracepoint IDs sufficiently long and descriptive names making the occurrence of collisions quite unlikely. The probability of collision is thus made sufficiently low so that the problem may, as a practical matter, be ignored. By requirement, the consequence of collisions will be a slight ambiguity in the trace streams; tracing will continue in spite of collisions and ambiguities. “The show must go on”. The *posix\_prog\_address* member of the **posix\_trace\_event\_info** structure is used to allow trace streams to be unambiguously interpreted, despite the fact that trace event types and trace event names need not be unique.

7564

7565

7566

7567

7568

7569

7570

7571

7572

7573

The *posix\_trace\_eventid\_open()* function is required to allow the instrumented third-party library to get a valid trace event type identifier for its trace event names. This operation is, somehow, an allocation, and the group was aware of proposing some deallocation mechanism which the instrumented application could use to recover the resources used by a trace event type identifier. This would have given the instrumented application the benefit of being capable of reusing a possible minimum set of trace event type identifiers, but also the inconvenience to have, possibly in the same trace stream, one trace event type identifier identifying two different trace event types. After some discussions the group decided to not define such a function which would make this API thicker for little benefit, the user having always the possibility of adding identification information in the data member of the trace event structure.

7574

7575

7576

7577

7578

The set of the trace event type identifiers the controlling process wants to filter out is initialized in the trace mechanism using the function *posix\_trace\_set\_filter()*, setting the arguments according to the definitions explained in *posix\_trace\_set\_filter()*. This operation can be done statically (when the trace is in the STOPPED state) or dynamically (when the trace is in the STARTED state). The preparation of the filter is normally done using the function defined in

7579 *posix\_trace\_eventtypelist\_getnext\_id()* and eventually the function  
7580 *posix\_trace\_eventtypelist\_rewind()* in order to know (before the recording) the list of the potential  
7581 set of trace event types that can be recorded. In the case of an active trace stream, this list may  
7582 not be exhaustive. Actually, the target process may not have yet called the function  
7583 *posix\_trace\_eventid\_open()*. But it is a common practice, for a controlling process, to prepare the  
7584 filtering of a future trace stream before its start. Therefore the user must have a way to get the  
7585 trace event type identifier corresponding to a well-known trace event name before its future  
7586 association by the pre-cited function. This is done by calling the *posix\_trace\_trid\_eventid\_open()*  
7587 function, given the trace stream identifier and the trace name, and described hereafter. Because  
7588 this trace event type identifier is associated with a trace stream identifier, where a unique  
7589 process has initialized two or more traces, the implementation is expected to return the same  
7590 trace event type identifier for successive calls to *posix\_trace\_trid\_eventid\_open()* with different  
7591 trace stream identifiers. The *posix\_trace\_eventid\_get\_name()* function is used by the controller  
7592 process to identify, by the name, the trace event type returned by a call to the  
7593 *posix\_trace\_eventtypelist\_getnext\_id()* function.

7594 Afterwards, the set of trace event types is constructed using the functions defined in  
7595 *posix\_trace\_eventset\_empty()*, *posix\_trace\_eventset\_fill()*, *posix\_trace\_eventset\_add()*, and  
7596 *posix\_trace\_eventset\_del()*.

7597 A set of functions is provided devoted to the manipulation of the trace event type identifier and  
7598 names for an active trace stream. All these functions require the trace stream identifier argument  
7599 as the first parameter. The opacity of the trace event type identifier implies that the user cannot  
7600 associate directly its well-known trace event name with the system associated trace event type  
7601 identifier.

7602 The *posix\_trace\_trid\_eventid\_open()* function allows the application to get the system trace event  
7603 type identifier back from the system, given its well-known trace event name. One possible use of  
7604 this function is to qualify a filter.

7605 The *posix\_trace\_eventid\_get\_name()* function allows the application to obtain a trace event name  
7606 given its trace event type identifier. One possible use of this function is to identify the type of a  
7607 trace event retrieved from the trace stream, and print it. The easiest way to implement this  
7608 requirement, is to use a single trace event type map for all the processes whose maps are  
7609 required to be identical. A more difficult way is to attempt to keep multiple maps identical at  
7610 every call to *posix\_trace\_eventid\_open()* and *posix\_trace\_trid\_eventid\_open()*.

#### 7611 *B.2.11.6 Rationale on Trace Events Type Filtering*

7612 The most basic rationale for runtime and pre-registration filtering (selection/rejection) of trace  
7613 event types is to prevent choking of the trace collection facility, and/or overloading of the  
7614 computer system. Any worthwhile trace facility can bring even the largest computer to its  
7615 knees. Otherwise, we would record everything, and filter after the fact; it would be much  
7616 simpler, but impractical.

7617 To achieve debugging, measurement, or whatever the purpose of tracing, the filtering of trace  
7618 event types is an important part of trace analysis. Due to the fact that the trace events are put  
7619 into a trace stream and probably logged afterwards into a file, different levels of filtering—that  
7620 is, rejection of trace event types—are possible.

7621 **Filtering of Trace Event Types Before Tracing**

7622 This function, represented by the `posix_trace_set_filter()` function in IEEE Std. 1003.1-200x (see  
 7623 `posix_trace_set_filter()`), selects, before or during tracing, the set of trace event types to be filtered  
 7624 out. It should be possible also (as OSF suggested in their ETAP trace specifications) to select the  
 7625 kernel trace event types to be traced in a system-wide fashion. These two functionalities are  
 7626 called the pre-filtering of trace event types.

7627 The restriction on the actual type used for the `trace_event_set_t` type is intended to guarantee  
 7628 that these objects can always be assigned, have their address taken, and be passed by value as  
 7629 parameters. It is not intended that this type be a structure including pointers to other data  
 7630 structures, as that could impact the portability of applications performing such operations. A  
 7631 reasonable implementation could be a structure containing an array of integer types.

7632 **Filtering of Trace Event Types at Runtime**

7633 Using this API, this functionality may be built, a privileged process or a privileged thread can  
 7634 get trace events from the trace stream of another process or thread, and thus specify the type of  
 7635 trace events to record into a file, using methods and interfaces out of the scope of  
 7636 IEEE Std. 1003.1-200x. This functionality, called inline filtering of trace event types, is used for  
 7637 runtime analysis of trace streams.

7638 **Post-Mortem Filtering of Trace Event Types**

7639 The word *post-mortem* is used here to indicate that some unanticipated situation occurs during  
 7640 execution that does not permit a pre or inline filtering of trace events and that it is necessary to  
 7641 record all trace event types, to have a chance to discover the problem afterwards. When the  
 7642 program stops, all the trace events recorded previously can be analyzed in order to find the  
 7643 solution. This functionality could be named the post-filtering of trace event types.

7644 **Discussions about Trace Event Type-Filtering**

7645 After long discussions with the parties involved in the process of defining the trace interface, it  
 7646 seems that the sensitivity to the filtering problem is different, but everybody agrees that the level  
 7647 of the overhead introduced during the tracing operation depends on the filtering method  
 7648 elected. If the time that it takes the trace event to be recorded can be neglected, the overhead  
 7649 introduced by the filtering process can be classified as follows:

7650 Pre-filtering      System and process/thread-level overhead

7651 Inline-filtering    Process/thread-level overhead

7652 Post-filtering     No overhead; done offline

7653 The pre-filtering could be named *critical realtime* filtering in the sense that the filtering of trace  
 7654 event type is manageable at the user level so the user can lower to a minimum the filtering  
 7655 overhead at some user selected level of priority for the inline filtering, or delay the filtering to  
 7656 after execution for the post-filtering. The counterpart of this solution is that the size of the trace  
 7657 stream must be sufficient to record all the trace events. The advantage of the pre-filtering is that  
 7658 the utilization of the trace stream is optimized.

7659 Only pre-filtering is defined by IEEE Std. 1003.1-200x. However, great care must be taken in  
 7660 specifying pre-filtering, so that it does not impose unacceptable overhead. Moreover, it is  
 7661 necessary to isolate all the functionality relative to the pre-filtering.

7662 The result of this rationale is to define a new option, the Trace Event Filter option, not  
 7663 necessarily implemented in small realtime systems, where system overhead is minimized to the  
 7664 extent possible.

7665 *B.2.11.7 Tracing, pthread API*

7666 The objective to be able to control tracing for individual threads may be in conflict with the  
 7667 efficiency expected in threads with a `contentionscope` attribute of  
 7668 `PTHREAD_SCOPE_PROCESS`. For these threads, context switches from one thread that has  
 7669 tracing enabled to another thread that has tracing disabled may require a kernel call to inform  
 7670 the kernel whether it has to trace system events executed by that thread or not. For this reason, it  
 7671 was proposed that the ability to enable or disable tracing for `PTHREAD_SCOPE_PROCESS`  
 7672 threads be made optional, through the introduction of a Trace Scope Process option. A trace  
 7673 implementation which did not implement the Trace Scope Process option would not honor the  
 7674 tracing-state attribute of a thread with `PTHREAD_SCOPE_PROCESS`; it would, however, honor  
 7675 the tracing-state attribute of a thread with `PTHREAD_SCOPE_SYSTEM`. This proposal was  
 7676 rejected as:

- 7677 1. Removing desired functionality (per-thread trace control)
- 7678 2. Introducing counter-intuitive behavior for the tracing-state attribute
- 7679 3. Mixing logically orthogonal ideas (thread scheduling and thread tracing)
- 7680 [Objective 4]

7681 Finally, to solve this complex issue, this API does not provide `pthread_gettracingstate()`,  
 7682 `pthread_settracingstate()`, `pthread_attr_gettracingstate()`, and `pthread_attr_settracingstate()`  
 7683 interfaces. These interfaces force the thread implementation to add to the weight of the thread  
 7684 and cause a revision of the threads libraries, just to support tracing. Worse yet,  
 7685 `posix_trace_userevent()` must always test this per-thread variable even in the common case where  
 7686 it is not used at all. Per-thread tracing is easy to implement using existing interfaces where  
 7687 necessary; see the following example.

7688 **Example**

```
7689 /* Caution. Error checks omitted */
7690 static pthread_key_t my_key;
7691 static trace_event_id_t my_event_id;
7692 static pthread_once_t my_once = PTHREAD_ONCE_INIT;

7693 void my_init(void)
7694 {
7695     (void) pthread_key_create(&my_key, NULL);
7696     (void) posix_trace_eventid_open("my", &my_event_id);
7697 }

7698 int get_trace_flag(void)
7699 {
7700     pthread_once(&my_once, my_init);
7701     return (pthread_getspecific(my_key) != NULL);
7702 }

7703 void set_trace_flag(int f)
7704 {
7705     pthread_once(&my_once, my_init);
7706     pthread_setspecific(my_key, f? &my_event_id: NULL);
7707 }

7708 fn()
7709 {
7710     if (get_trace_flag())
```

```
7711         posix_trace_event(my_event_id, ...)
7712     }
```

7713 The above example does not implement third-party state setting, but it is also implementable  
7714 with some more work, yet the extra functionality is rarely needed.

7715 Lastly, per-thread tracing works poorly for threads with PTHREAD\_SCOPE\_PROCESS  
7716 contention scope. These “library” threads have minimal interaction with the kernel and would  
7717 have to explicitly set the attributes whenever they are context switched to a new kernel thread in  
7718 order to trace system events. Such state was explicitly avoided in POSIX threads to keep  
7719 PTHREAD\_SCOPE\_PROCESS threads lightweight.

7720 The reason that keeping PTHREAD\_SCOPE\_PROCESS threads lightweight is important is that  
7721 such threads can be used not just for simple multi-processors but also for coroutine style  
7722 programming (such as discrete event simulation) without inventing a new threads paradigm.  
7723 Adding extra runtime cost to thread context switches will make using POSIX threads less  
7724 attractive in these situations.

#### 7725 *B.2.11.8 Rationale on Triggering*

7726 The ability to start or stop tracing based on the occurrence of specific trace event types has been  
7727 proposed as a parallel to similar functionality appearing in logic analyzers. Such triggering, in  
7728 order to be very useful, should be based not only on the trace event type, but on trace event-  
7729 specific data, including tests of user-specified fields for matching or threshold values.

7730 Such a facility is unnecessary where the buffering of the stream is not a constraint, since such  
7731 checks can be performed offline during post-mortem analysis.

7732 For example, a large system could incorporate a daemon utility to collect the trace records from  
7733 memory buffers and spool them to secondary storage for later analysis. In the instances where  
7734 resources are truly limited, such as embedded applications, the application incorporation of  
7735 application code to test the circumstances of a trace event and call the trace point only if needed  
7736 is usually straightforward.

7737 For performance reasons, the *posix\_trace\_event()* function should be implemented using a macro,  
7738 so if the trace is inactive, the trace event point calls are latent code and must cost no more than a  
7739 scalar test.

7740 The API proposed in IEEE Std. 1003.1-200x does not include any triggering functionality.

#### 7741 *B.2.11.9 Rationale on Timestamp Clock*

7742 It has been suggested that the tracing mechanism should include the possibility of specifying the  
7743 clock to be used in timestamping the trace events. When application trace events must be  
7744 correlated to remote trace events, such a facility could provide a global time reference not  
7745 available from a local clock. Further, the application may be driven by timers based on a clock  
7746 different from that used for the timestamp, and the correlation of the trace to those untraced  
7747 timer activities could be an important part of the analysis of the application.

7748 However, the tracing mechanism needs to be fast and just the provision of such an option can  
7749 materially affect its performance. Leaving aside the performance costs of reading some clocks,  
7750 this notion is also ill-defined when kernel trace events are to be traced by two applications  
7751 making use of different tracing clocks. This can even happen within a single application where  
7752 different parts of the application are served by different clocks. Another complication can occur  
7753 when a clock is maintained strictly at the user level and is unavailable at the kernel level.

7754 It is felt that the benefits of a selectable trace clock do not match its costs. Applications that wish  
7755 to correlate clocks other than the default tracing clock can include trace events with sample

7756 values of those other clocks, allowing correlation of timestamps from the various independent  
7757 clocks. In any case, such a technique would be required when applications are sensitive to  
7758 multiple clocks.

#### 7759 *B.2.11.10 Rationale on Different Overrun Conditions*

7760 The analysis of the dynamic behavior of the trace mechanism shows that different overrun  
7761 conditions may occur. The API must provide a means to manage such conditions in a portable  
7762 way.

#### 7763 **Overrun in Trace Streams Initialized with POSIX\_TRACE\_LOOP Policy**

7764 In this case, the user of the trace mechanism is interested in using the trace stream with  
7765 POSIX\_TRACE\_LOOP policy to record trace events continuously, but ideally without losing any  
7766 trace events. The online analyzer process must get the trace events at a mean speed equivalent to  
7767 the recording speed. Should the trace stream become full, a trace stream overrun occurs. This  
7768 condition is detected by getting the status of the active trace stream (function  
7769 *posix\_trace\_get\_status()*) and looking at the member *posix\_stream\_overrun\_status* of the read  
7770 **posix\_stream\_status** structure. In addition, two predefined trace event types are defined:

- 7771 1. The beginning of a trace overflow, to locate the beginning of an overflow when reading a  
7772 trace stream
- 7773 2. The end of a trace overflow, to locate the end of an overflow, when reading a trace stream

7774 As a timestamp is associated with these predefined trace events, it is possible to know the  
7775 duration of the overflow.

#### 7776 **Overrun in Dumping Trace Streams into Trace Logs**

7777 The user lets the trace mechanism dump the trace stream initialized with  
7778 POSIX\_TRACE\_FLUSH policy automatically into a trace log. If the dump operation is slower  
7779 than the recording of trace events, the trace stream can overrun. This condition is detected by  
7780 getting the status of the active trace stream (function *posix\_trace\_get\_status()*) and looking at the  
7781 member *posix\_log\_overrun\_status* of the read **posix\_stream\_status** structure. This overrun  
7782 indicates that the trace mechanism is not able to operate in this mode at this speed. It is the  
7783 responsibility of the user to modify one of the trace parameters (the stream size or the trace  
7784 event type filter, for instance) to avoid such overrun conditions, if overruns are to be prevented.  
7785 The same already predefined trace event types (see **Overrun in Trace Streams Initialized with**  
7786 **POSIX\_TRACE\_LOOP Policy**) are used to detect and to know the duration of an overflow.

#### 7787 **Reading an Active Trace Stream**

7788 Although this trace API allows one to read an active trace stream with log while it is tracing, this  
7789 feature can lead to false overflow origin interpretation: the trace log or the reader of the trace  
7790 stream. Reading from an active trace stream with log is thus non-portable, and has been left  
7791 unspecified.

7792 **B.2.12 Data Types**

7793 The requirement that additional types defined in this section end in “\_t” was prompted by the  
 7794 problem of name space pollution. It is difficult to define a type (where that type is not one  
 7795 defined by IEEE Std. 1003.1-200x) in one header file and use it in another without adding  
 7796 symbols to the name space of the program. To allow implementors to provide their own types,  
 7797 all conforming applications are required to avoid symbols ending in “\_t”, which permits the  
 7798 implementor to provide additional types. Because a major use of types is in the definition of  
 7799 structure members, which can (and in many cases must) be added to the structures defined in  
 7800 IEEE Std. 1003.1-200x, the need for additional types is compelling.

7801 The types, such as **ushort** and **ulong**, which are in common usage, are not defined in  
 7802 IEEE Std. 1003.1-200x (although **ushort\_t** would be permitted as an extension). They can be  
 7803 added to `<sys/types.h>` using a feature test macro (see Section B.2.2.1 (on page 3384)). A  
 7804 suggested symbol for these is `_SYSIII`. Similarly, the types like **u\_short** would probably be best  
 7805 controlled by `_BSD`.

7806 Some of these symbols may appear in other headers; see Section B.2.2.2 (on page 3384).

7807 **dev\_t** This type may be made large enough to accommodate host-locality considerations  
 7808 of networked systems.

7809 This type must be arithmetic. Earlier proposals allowed this to be non-arithmetic  
 7810 (such as a structure) and provided a `samefile()` function for comparison.

7811 **gid\_t** Some implementations had separated **gid\_t** from **uid\_t** before POSIX.1 was  
 7812 completed. It would be difficult for them to coalesce them when it was  
 7813 unnecessary. Additionally, it is quite possible that user IDs might be different than  
 7814 group IDs because the user ID might wish to span a heterogeneous network,  
 7815 where the group ID might not.

7816 For current implementations, the cost of having a separate **gid\_t** will be only  
 7817 lexical.

7818 **mode\_t** This type was chosen so that implementations could choose the appropriate  
 7819 integral type, and for compatibility with the ISO C standard. 4.3 BSD uses  
 7820 **unsigned short** and the SVID uses **ushort**, which is the same. Historically, only the  
 7821 low-order sixteen bits are significant.

7822 **nlink\_t** This type was introduced in place of **short** for `st_nlink` (see the `<sys/stat.h>` header)  
 7823 in response to an objection that **short** was too small.

7824 **off\_t** This type is used only in `lseek()`, `fcntl()`, and `<sys/stat.h>`. Many implementations  
 7825 would have difficulties if it were defined as anything other than **long**. Requiring  
 7826 an integral type limits the capabilities of `lseek()` to four gigabytes. The ISO C  
 7827 standard supplies routines that use larger types; see `fgetpos()` and `fsetpos()`. XSI-  
 7828 conformant systems provide the `fseeko()` and `lseeko()` functions that use larger  
 7829 types.

7830 **pid\_t** The inclusion of this symbol was controversial because it is tied to the issue of the  
 7831 representation of a process ID as a number. From the point of view of a portable  
 7832 application, process IDs should be “magic cookies”<sup>1</sup> that are produced by calls

7833 \_\_\_\_\_

7834 1. An historical term meaning: “An opaque object, or token, of determinate size, whose significance is known only to the entity  
 7835 which created it. An entity receiving such a token from the generating entity may only make such use of the ‘cookie’ as is defined  
 7836 and permitted by the supplying entity.”



7837 such as *fork()*, used by calls such as *waitpid()* or *kill()*, and not otherwise analyzed  
 7838 (except that the sign is used as a flag for certain operations).

7839 The concept of a {PID\_MAX} value interacted with this in early proposals. Treating  
 7840 process IDs as an opaque type both removes the requirement for {PID\_MAX} and  
 7841 allows systems to be more flexible in providing process IDs that span a large range  
 7842 of values, or a small one.

7843 Since the values in **uid\_t**, **gid\_t**, and **pid\_t** will be numbers generally, and  
 7844 potentially both large in magnitude and sparse, applications that are based on  
 7845 arrays of objects of this type are unlikely to be fully portable in any case. Solutions  
 7846 that treat them as magic cookies will be portable.

7847 {CHILD\_MAX} precludes the possibility of a “toy implementation”, where there  
 7848 would only be one process.

7849 **ssize\_t** This is intended to be a signed analog of **size\_t**. The wording is such that an  
 7850 implementation may either choose to use a longer type or simply to use the signed  
 7851 version of the type that underlies **size\_t**. All functions that return **ssize\_t** (*read()*  
 7852 and *write()*) describe as “implementation-defined” the result of an input exceeding  
 7853 {SSIZE\_MAX}. It is recognized that some implementations might have **ints** that  
 7854 are smaller than **size\_t**. A portable application would be constrained not to  
 7855 perform I/O in pieces larger than {SSIZE\_MAX}, but a portable application using  
 7856 extensions would be able to use the full range if the implementation provided an  
 7857 extended range, while still having a single type-compatible interface.

7858 The symbols **size\_t** and **ssize\_t** are also required in **<unistd.h>** to minimize the  
 7859 changes needed for calls to *read()* and *write()*. Implementors are reminded that it  
 7860 must be possible to include both **<sys/types.h>** and **<unistd.h>** in the same  
 7861 program (in either order) without error.

7862 **uid\_t** Before the addition of this type, the data types used to represent these values  
 7863 varied throughout early proposals. The **<sys/stat.h>** header defined these values as  
 7864 type **short**, the **<passwd.h>** file (now **<pwd.h>** and **<grp.h>**) used an **int**, and  
 7865 *getuid()* returned an **int**. In response to a strong objection to the inconsistent  
 7866 definitions, all the types to were switched to **uid\_t**.

7867 In practice, those historical implementations that use varying types of this sort can  
 7868 typedef **uid\_t** to **short** with no serious consequences.

7869 The problem associated with this change concerns object compatibility after  
 7870 structure size changes. Since most implementations will define **uid\_t** as a short, the  
 7871 only substantive change will be a reduction in the size of the **passwd** structure.  
 7872 Consequently, implementations with an overriding concern for object  
 7873 compatibility can pad the structure back to its current size. For that reason, this  
 7874 problem was not considered critical enough to warrant the addition of a separate  
 7875 type to POSIX.1.

7876 The types **uid\_t** and **gid\_t** are magic cookies. There is no {UID\_MAX} defined by  
 7877 POSIX.1, and no structure imposed on **uid\_t** and **gid\_t** other than that they be  
 7878 positive arithmetic types. (In fact, they could be **unsigned char**.) There is no  
 7879 maximum or minimum specified for the number of distinct user or group IDs.

7880 **B.3 System Interfaces**

7881 See the RATIONALE sections on the individual reference pages.

7882 **B.3.1 Examples for Spawn**7883 The following long examples are provided in the Rationale (Informative) volume of  
7884 IEEE Std. 1003.1-200x as a supplement to the reference page for *spawn()*.7885 **Example Library Implementation of Spawn**7886 The *posix\_spawn()* or *posix\_spawnnp()* functions provide the following:

- 7887 • Simply start a process executing a process image. This is the simplest application for process  
7888 creation, and it may cover most executions of POSIX *fork()*.
- 7889 • Support I/O redirection, including pipes.
- 7890 • Run the child under a user and group ID in the domain of the parent.
- 7891 • Run the child at any priority in the domain of the parent.

7892 The *posix\_spawn()* or *posix\_spawnnp()* functions do not cover every possible use of the *fork()*  
7893 function, but they do span the common applications: typical use by a shell and a login utility.7894 The price for an application is that before it calls *posix\_spawn()* or *posix\_spawnnp()*, the parent  
7895 must adjust to a state that *posix\_spawn()* or *posix\_spawnnp()* can map to the desired state for the  
7896 child. Environment changes require the parent to save some of its state and restore it afterwards.  
7897 The effective behavior of a successful invocation of *posix\_spawn()* is as if the operation were  
7898 implemented with POSIX operations as follows:

```

7899 #include <sys/types.h>
7900 #include <stdlib.h>
7901 #include <stdio.h>
7902 #include <unistd.h>
7903 #include <sched.h>
7904 #include <fcntl.h>
7905 #include <signal.h>
7906 #include <errno.h>
7907 #include <string.h>
7908 #include <signal.h>

7909 /* #include <spawn.h>*/
7910 /*****
7911  /* Things that could be defined in spawn.h */
7912  *****/
7913 typedef struct
7914 {
7915     short posix_attr_flags;
7916 #define POSIX_SPAWN_SETPGROUP          0x1
7917 #define POSIX_SPAWN_SETSIGMASK        0x2
7918 #define POSIX_SPAWN_SETSIGDEF         0x4
7919 #define POSIX_SPAWN_SETSCHEDULER      0x8
7920 #define POSIX_SPAWN_SETSCHEDPARAM     0x10
7921 #define POSIX_SPAWN_RESETIDS          0x20
7922     pid_t posix_attr_pgroup;
7923     sigset_t posix_attr_sigmask;
7924     sigset_t posix_attr_sigdefault;

```

```

7925     int posix_attr_schedpolicy;
7926     struct sched_param posix_attr_schedparam;
7927     } posix_spawnattr_t;
7928
7928     typedef char *posix_spawn_file_actions_t;
7929
7929     int posix_spawn_file_actions_init(
7930         posix_spawn_file_actions_t *file_actions);
7931     int posix_spawn_file_actions_destroy(
7932         posix_spawn_file_actions_t *file_actions);
7933     int posix_spawn_file_actions_addclose(
7934         posix_spawn_file_actions_t *file_actions, int fildes);
7935     int posix_spawn_file_actions_adddup2(
7936         posix_spawn_file_actions_t *file_actions, int fildes,
7937         int newfildes);
7938     int posix_spawn_file_actions_addopen(
7939         posix_spawn_file_actions_t *file_actions, int fildes,
7940         const char *path, int oflag, mode_t mode);
7941     int posix_spawnattr_init(posix_spawnattr_t *attr);
7942     int posix_spawnattr_destroy(posix_spawnattr_t *attr);
7943     int posix_spawnattr_getflags(const posix_spawnattr_t *attr, short *lags);
7944     int posix_spawnattr_setflags(posix_spawnattr_t *attr, short flags);
7945     int posix_spawnattr_getpgroup(const posix_spawnattr_t *attr,
7946         pid_t *pgroup);
7947     int posix_spawnattr_setpgroup(posix_spawnattr_t *attr, pid_t pgroup);
7948     int posix_spawnattr_getschedpolicy(const posix_spawnattr_t *attr,
7949         int *schedpolicy);
7950     int posix_spawnattr_setschedpolicy(posix_spawnattr_t *attr,
7951         int schedpolicy);
7952     int posix_spawnattr_getschedparam(const posix_spawnattr_t *attr,
7953         struct sched_param *schedparam);
7954     int posix_spawnattr_setschedparam(posix_spawnattr_t *attr,
7955         const struct sched_param *schedparam);
7956     int posix_spawnattr_getsigmask(const posix_spawnattr_t *attr,
7957         sigset_t *sigmask);
7958     int posix_spawnattr_setsigmask(posix_spawnattr_t *attr,
7959         const sigset_t *sigmask);
7960     int posix_spawnattr_getdefault(const posix_spawnattr_t *attr,
7961         sigset_t *sigdefault);
7962     int posix_spawnattr_setsigdefault(posix_spawnattr_t *attr,
7963         const sigset_t *sigdefault);
7964     int posix_spawn(pid_t *pid, const char *path,
7965         const posix_spawn_file_actions_t *file_actions,
7966         const posix_spawnattr_t *attrp, char * const argv[],
7967         char * const envp[]);
7968     int posix_spawnnp(pid_t *pid, const char *file,
7969         const posix_spawn_file_actions_t *file_actions,
7970         const posix_spawnattr_t *attrp, char * const argv[],
7971         char * const envp[]);
7972
7972     /*****
7973     /* Example posix_spawn() library routine */
7974     /*****
7975     int posix_spawn(pid_t *pid,

```

```
7976     const char *path,
7977     const posix_spawn_file_actions_t *file_actions,
7978     const posix_spawnattr_t *attrp,
7979     char * const argv[],
7980     char * const envp[])
7981 {
7982     /* Create process */
7983     if((*pid=fork()) == (pid_t)0)
7984     {
7985         /* This is the child process */
7986         /* Worry about process group */
7987         if(attrp->posix_attr_flags & POSIX_SPAWN_SETPGROUP)
7988         {
7989             /* Override inherited process group */
7990             if(setpgid(0, attrp->posix_attr_pgroup) != 0)
7991             {
7992                 /* Failed */
7993                 exit(127);
7994             }
7995         }
7996
7997         /* Worry about process signal mask */
7998         if(attrp->posix_attr_flags & POSIX_SPAWN_SETSIGMASK)
7999         {
8000             /* Set the signal mask (can't fail) */
8001             sigprocmask(SIG_SETMASK , &attrp->posix_attr_sigmask,
8002             NULL);
8003         }
8004
8005         /* Worry about resetting effective user and group IDs */
8006         if(attrp->posix_attr_flags & POSIX_SPAWN_RESETEIDS)
8007         {
8008             /* None of these can fail for this case. */
8009             setuid(getuid());
8010             setgid(getgid());
8011         }
8012
8013         /* Worry about defaulted signals */
8014         if(attrp->posix_attr_flags & POSIX_SPAWN_SETSIGDEF)
8015         {
8016             struct sigaction deflt;
8017             sigset_t all_signals;
8018
8019             int s;
8020
8021             /* Construct default signal action */
8022             deflt.sa_handler = SIG_DFL;
8023             deflt.sa_flags = 0;
8024
8025             /* Construct the set of all signals */
8026             sigfillset(&all_signals);
8027
8028             /* Loop for all signals */
8029             for(s=0; sigismember(&all_signals,s); s++)
8030             {
8031                 /* Signal to be defaulted? */
```

```

8025         if(sigismember(&attrp->posix_attr_sigdefault,s))
8026             {
8027                 /* Yes; default this signal */
8028                 if(sigaction(s, &deflt, NULL) == -1)
8029                     {
8030                         /* Failed */
8031                         exit(127);
8032                     }
8033             }
8034     }
8035 }

8036 /* Worry about the fds if we are to map them */
8037 if(file_actions != NULL)
8038     {
8039     /* Loop for all actions in object file_actions */
8040     /*(implementation dives beneath abstraction)*/
8041     char *p = *file_actions;
8042     while(*p != ' ')
8043         {
8044             if(strncmp(p,"close(",6) == 0)
8045                 {
8046                     int fd;
8047                     if(sscanf(p+6,"%d",&fd) != 1)
8048                         {
8049                             exit(127);
8050                         }
8051                     if(close(fd) == -1) exit(127);
8052                 }
8053             else if(strncmp(p,"dup2(",5) == 0)
8054                 {
8055                     int fd,newfd;
8056                     if(sscanf(p+5,"%d,%d",&fd,&newfd) != 2)
8057                         {
8058                             exit(127);
8059                         }
8060                     if(dup2(fd, newfd) == -1) exit(127);
8061                 }
8062             else if(strncmp(p,"open(",5) == 0)
8063                 {
8064                     int fd,oflag;
8065                     mode_t mode;
8066                     int tempfd;
8067                     char path[1000]; /* Should be dynamic */
8068                     char *q;
8069                     if(sscanf(p+5,"%d",&fd) != 1)
8070                         {
8071                             exit(127);
8072                         }
8073                     p = strchr(p, ' ') + 1;
8074                     q = strchr(p, '*');
8075                     if(q == NULL) exit(127);
8076                     strncpy(path, p, q-p);

```

```

8077         path[q-p] = ' ';
8078     if(sscanf(q+1,"%o,%o",&oflag,&mode)!=2)
8079     {
8080         exit(127);
8081     }
8082     if(close(fd) == -1)
8083     {
8084         if(errno != EBADF) exit(127);
8085     }
8086     tempfd = open(path, oflag, mode);
8087     if(tempfd == -1) exit(127);
8088     if(tempfd != fd)
8089     {
8090         if(dup2(tempfd,fd) == -1)
8091         {
8092             exit(127);
8093         }
8094         if(close(tempfd) == -1)
8095         {
8096             exit(127);
8097         }
8098     }
8099     }
8100     else
8101     {
8102         exit(127);
8103     }
8104     p = strchr(p, ' ') + 1;
8105 }
8106 }

8107 /* Worry about setting new scheduling policy and parameters */
8108 if(attrp->posix_attr_flags & POSIX_SPAWN_SETSCHEDULER)
8109 {
8110     if(sched_setscheduler(0, attrp->posix_attr_schedpolicy,
8111         &attrp->posix_attr_schedparam) == -1)
8112     {
8113         exit(127);
8114     }
8115 }

8116 /* Worry about setting only new scheduling parameters */
8117 if(attrp->posix_attr_flags & POSIX_SPAWN_SETSCHEDPARAM)
8118 {
8119     if(sched_setparam(0, &attrp->posix_attr_schedparam)==-1)
8120     {
8121         exit(127);
8122     }
8123 }

8124 /* Now execute the program at path */
8125 /* Any fd that still has FD_CLOEXEC set will be closed */
8126 execve(path, argv, envp);
8127 exit(127); /* exec failed */

```

```

8128         }
8129         else
8130         {
8131             /* This is the parent (calling) process */
8132             if(*pid == (pid_t)-1) return errno;
8133             return 0;
8134         }
8135     }

8136     /*****
8137     /* Here is a crude but effective implementation of the */
8138     /* file action object operators which store actions as */
8139     /* concatenated token separated strings.          */
8140     /*****
8141     /* Create object with no actions. */
8142     int posix_spawn_file_actions_init(
8143         posix_spawn_file_actions_t *file_actions)
8144     {
8145         *file_actions = malloc(sizeof(char));
8146         if(*file_actions == NULL) return ENOMEM;
8147         strcpy(*file_actions, "");
8148         return 0;
8149     }

8150     /* Free object storage and make invalid. */
8151     int posix_spawn_file_actions_destroy(
8152         posix_spawn_file_actions_t *file_actions)
8153     {
8154         free(*file_actions);
8155         *file_actions = NULL;
8156         return 0;
8157     }

8158     /* Add a new action string to object. */
8159     static int add_to_file_actions(
8160         posix_spawn_file_actions_t *file_actions,
8161         char *new_action)
8162     {
8163         *file_actions = realloc
8164             (*file_actions, strlen(*file_actions)+strlen(new_action)+1);
8165         if(*file_actions == NULL) return ENOMEM;
8166         strcat(*file_actions, new_action);
8167         return 0;
8168     }

8169     /* Add a close action to object. */
8170     int posix_spawn_file_actions_addclose(
8171         posix_spawn_file_actions_t *file_actions, int fildes)
8172     {
8173         char temp[100];
8174         sprintf(temp, "close(%d)", fildes);
8175         return add_to_file_actions(file_actions, temp);
8176     }

```

```

8177     /* Add a dup2 action to object. */
8178     int posix_spawn_file_actions_adddup2(
8179         posix_spawn_file_actions_t *file_actions, int fildes,
8180         int newfildes)
8181     {
8182         char temp[100];
8183         sprintf(temp, "dup2(%d,%d)", fildes, newfildes);
8184         return add_to_file_actions(file_actions, temp);
8185     }

8186     /* Add an open action to object. */
8187     int posix_spawn_file_actions_addopen(
8188         posix_spawn_file_actions_t *file_actions, int fildes,
8189         const char *path, int oflag, mode_t mode)
8190     {
8191         char temp[100];
8192         sprintf(temp, "open(%d,%s*%o,%o)", fildes, path, oflag, mode);
8193         return add_to_file_actions(file_actions, temp);
8194     }

8195     /*****
8196     /* Here is a crude but effective implementation of the */
8197     /* spawn attributes object functions which manipulate */
8198     /* the individual attributes. */
8199     /*****
8200     /* Initialize object with default values. */
8201     int posix_spawnattr_init(
8202         posix_spawnattr_t *attr)
8203     {
8204         attr->posix_attr_flags=0;
8205         attr->posix_attr_pgroup=0;
8206         /* Default value of signal mask is the parent's signal mask; */
8207         /* other values are also allowed */
8208         sigprocmask(0,NULL,&attr->posix_attr_sigmask);
8209         sigemptyset(&attr->posix_attr_sigdefault);
8210         /* Default values of scheduling attr inherited from the parent; */
8211         /* other values are also allowed */
8212         attr->posix_attr_schedpolicy=sched_getscheduler(0);
8213         sched_getparam(0,&attr->posix_attr_schedparam);
8214         return 0;
8215     }

8216     int posix_spawnattr_destroy(posix_spawnattr_t *attr)
8217     {
8218         /* No action needed */
8219         return 0;
8220     }

8221     int posix_spawnattr_getflags(const posix_spawnattr_t *attr,
8222         short *flags)
8223     {
8224         *flags=attr->posix_attr_flags;
8225         return 0;
8226     }

```



```
8227     int posix_spawnattr_setflags(posix_spawnattr_t *attr, short flags)
8228     {
8229         attr->posix_attr_flags=flags;
8230         return 0;
8231     }

8232     int posix_spawnattr_getpgroup(const posix_spawnattr_t *attr,
8233         pid_t *pgroup)
8234     {
8235         *pgroup=attr->posix_attr_pgroup;
8236         return 0;
8237     }

8238     int posix_spawnattr_setpgroup(posix_spawnattr_t *attr, pid_t pgroup)
8239     {
8240         attr->posix_attr_pgroup=pgroup;
8241         return 0;
8242     }

8243     int posix_spawnattr_getschedpolicy(const posix_spawnattr_t *attr,
8244         int *schedpolicy)
8245     {
8246         *schedpolicy=attr->posix_attr_schedpolicy;
8247         return 0;
8248     }

8249     int posix_spawnattr_setschedpolicy(posix_spawnattr_t *attr,
8250         int schedpolicy)
8251     {
8252         attr->posix_attr_schedpolicy=schedpolicy;
8253         return 0;
8254     }

8255     int posix_spawnattr_getschedparam(const posix_spawnattr_t *attr,
8256         struct sched_param *schedparam)
8257     {
8258         *schedparam=attr->posix_attr_schedparam;
8259         return 0;
8260     }

8261     int posix_spawnattr_setschedparam(posix_spawnattr_t *attr,
8262         const struct sched_param *schedparam)
8263     {
8264         attr->posix_attr_schedparam=*schedparam;
8265         return 0;
8266     }

8267     int posix_spawnattr_getsigmask(const posix_spawnattr_t *attr,
8268         sigset_t *sigmask)
8269     {
8270         *sigmask=attr->posix_attr_sigmask;
8271         return 0;
8272     }

8273     int posix_spawnattr_setsigmask(posix_spawnattr_t *attr,
8274         const sigset_t *sigmask)
```

```

8275     {
8276     attr->posix_attr_sigmask=*sigmask;
8277     return 0;
8278     }

8279     int posix_spawnattr_getsigdefault(const posix_spawnattr_t *attr,
8280     sigset_t *sigdefault)
8281     {
8282     *sigdefault=attr->posix_attr_sigdefault;
8283     return 0;
8284     }

8285     int posix_spawnattr_setsigdefault(posix_spawnattr_t *attr,
8286     const sigset_t *sigdefault)
8287     {
8288     attr->posix_attr_sigdefault=*sigdefault;
8289     return 0;
8290     }

```

### 8291 **I/O Redirection with Spawn**

8292 I/O redirection with *posix\_spawn()* or *posix\_spawnp()* is accomplished by crafting a *file\_actions*  
8293 argument to effect the desired redirection. Such a redirection follows the general outline of the  
8294 following example:

```

8295     /* To redirect new standard output (fd 1) to a file, */
8296     /* and redirect new standard input (fd 0) from my fd socket_pair[1], */
8297     /* and close my fd socket_pair[0] in the new process. */
8298     posix_spawn_file_actions_t file_actions;
8299     posix_spawn_file_actions_init(&file_actions);
8300     posix_spawn_file_actions_addopen(&file_actions, 1, "newout", ...);
8301     posix_spawn_file_actions_dup2(&file_actions, socket_pair[1], 0);
8302     posix_spawn_file_actions_close(&file_actions, socket_pair[0]);
8303     posix_spawn_file_actions_close(&file_actions, socket_pair[1]);
8304     posix_spawn(..., &file_actions, ...);
8305     posix_spawn_file_actions_destroy(&file_actions);

```

### 8306 **Spawning a Process Under a New User ID**

8307 Spawning a process under a new user ID follows the outline shown in the following example:

```

8308     Save = getuid();
8309     setuid(newid);
8310     posix_spawn(...);
8311     setuid(Save);

```

8312 / *Rationale (Informative)*

8313 **Part C:**

8314 **Shell and Utilities**

8315 *The Open Group*



# Rationale for Shell and Utilities

8316

## 8317 **C.1 Introduction**

### 8318 **C.1.1 Scope**

8319 Refer to Section A.1.1 (on page 3311).

### 8320 **C.1.2 Conformance**

8321 Refer to Section A.2 (on page 3317).

### 8322 **C.1.3 Normative References**

8323 There is no additional rationale provided for this section.

### 8324 **C.1.4 Changes from Issue 4**

8325 The change history is provided as an informative section, to track changes from previous issues  
8326 of IEEE Std. 1003.1-200x that comprised earlier versions of the Single UNIX Specification.

#### 8327 *C.1.4.1 Changes from Issue 4 to Issue 4, Version 2*

8328 There is no additional rationale provided for this section.

#### 8329 *C.1.4.2 Changes from Issue 4, Version 2 to Issue 5*

8330 There is no additional rationale provided for this section.

#### 8331 *C.1.4.3 Changes from Issue 5 to Issue 6*

8332 There is no additional rationale provided for this section.

### 8333 **C.1.5 Terminology**

8334 Refer to Section A.1.4 (on page 3313).

### 8335 **C.1.6 Definitions**

8336 Refer to Section A.3 (on page 3321).

### 8337 **C.1.7 Relationship to Other Documents**

#### 8338 *C.1.7.1 System Interfaces*

8339 It has been pointed out that the Shell and Utilities volume of IEEE Std. 1003.1-200x assumes that  
8340 a great deal of functionality from the System Interfaces volume of IEEE Std. 1003.1-200x is  
8341 present, but never states exactly how much (and strictly does not need to since both are  
8342 mandated on a conforming system). This section is an attempt to clarify the assumptions.

**8343 C.1.8 Portability**

8344 Refer to Section A.1.5 (on page 3315).

**8345 C.1.8.1 Codes**

8346 Refer to Section A.1.5.1 (on page 3315).

**8347 C.1.9 Utility Limits**

8348 This section grew out of an idea that originated with the original POSIX.1, in the tables of system  
8349 limits for the *sysconf()* and *pathconf()* functions. The idea being that a conforming application  
8350 can be written to use the most restrictive values that a minimal system can provide, but it should  
8351 not have to. The values provided represent compromises so that some vendors can use  
8352 historically limited versions of UNIX system utilities. They are the highest values that a strictly  
8353 conforming application can assume, given no other information.

8354 However, by using the *getconf* utility or the *sysconf()* function, the elegant application can be  
8355 tailored to more liberal values on some of the specific instances of specific implementations.

8356 There is no explicitly stated requirement that an implementation provide finite limits for any of  
8357 these numeric values; the implementation is free to provide essentially unbounded capabilities  
8358 (where it makes sense), stopping only at reasonable points such as {ULONG\_MAX} (from the  
8359 ISO C standard). Therefore, applications desiring to tailor themselves to the values on a  
8360 particular implementation need to be ready for possibly huge values; it may not be a good idea  
8361 to allocate blindly a buffer for an input line based on the value of {LINE\_MAX}, for instance.  
8362 However, unlike the System Interfaces volume of IEEE Std. 1003.1-200x, there is no set of limits  
8363 that return a special indication meaning “unbounded”. The implementation should always  
8364 return an actual number, even if the number is very large.

8365 The statement:

8366 “It is not guaranteed that the application ...”

8367 is an indication that many of these limits are designed to ensure that implementors design their  
8368 utilities without arbitrary constraints related to unimaginative programming. There are certainly  
8369 conditions under which combinations of options can cause failures that would not render an  
8370 implementation non-conforming. For example, {EXPR\_NEST\_MAX} and {ARG\_MAX} could  
8371 collide when expressions are large; combinations of {BC\_SCALE\_MAX} and {BC\_DIM\_MAX}  
8372 could exceed virtual memory.

8373 In the Shell and Utilities volume of IEEE Std. 1003.1-200x, the notion of a limit being guaranteed  
8374 for the process lifetime, as it is in the System Interfaces volume of IEEE Std. 1003.1-200x, is not as  
8375 useful to a shell script. The *getconf* utility is probably a process itself, so the guarantee would be  
8376 without value. Therefore, the Shell and Utilities volume of IEEE Std. 1003.1-200x requires the  
8377 guarantee to be for the session lifetime. This will mean that many vendors will either return very  
8378 conservative values or possibly implement *getconf* as a built-in.

8379 It may seem confusing to have limits that apply only to a single utility grouped into one global  
8380 section. However, the alternative, which would be to disperse them out into their utility  
8381 description sections, would cause great difficulty when *sysconf()* and *getconf* were described.  
8382 Therefore, the standard developers chose the global approach.

8383 Each language binding could provide symbol names that are slightly different than are shown  
8384 here. For example, the C-Language Binding option adds a leading underscore to the symbols as a  
8385 prefix.

8386 The following comments describe selection criteria for the symbols and their values:

8387 {ARG\_MAX}

8388 This is defined by the System Interfaces volume of IEEE Std. 1003.1-200x. Unfortunately, it

8389 is very difficult for a portable application to deal with this value, as it does not know how

8390 much of its argument space is being consumed by the environment variables of the user.

8391 {BC\_BASE\_MAX}

8392 {BC\_DIM\_MAX}

8393 {BC\_SCALE\_MAX}

8394 These were originally one value, {BC\_SCALE\_MAX}, but it was unreasonable to link all

8395 three concepts into one limit.

8396 {CHILD\_MAX}

8397 This is defined by the System Interfaces volume of IEEE Std. 1003.1-200x.

8398 {COLL\_WEIGHTS\_MAX}

8399 The weights assigned to **order** can be considered as “passes” through the collation

8400 algorithm.

8401 {EXPR\_NEST\_MAX}

8402 The value for expression nesting was borrowed from the ISO C standard.

8403 {LINE\_MAX}

8404 This is a global limit that affects all utilities, unless otherwise noted. The {MAX\_CANON}

8405 value from the System Interfaces volume of IEEE Std. 1003.1-200x may further limit input

8406 lines from terminals. The {LINE\_MAX} value was the subject of much debate and is a

8407 compromise between those who wished to have unlimited lines and those who understood

8408 that many historical utilities were written with fixed buffers. Frequently, utility writers

8409 selected the UNIX system constant {BUFSIZ} to allocate these buffers; therefore, some

8410 utilities were limited to 512 bytes for I/O lines, while others achieved 4 096 bytes or greater.

8411 It should be noted that {LINE\_MAX} applies only to input line length; there is no

8412 requirement in IEEE Std. 1003.1-200x that limits the length of output lines. Utilities such as

8413 *awk*, *sed*, and *paste* could theoretically construct lines longer than any of the input lines they

8414 received, depending on the options used or the instructions from the application. They are

8415 not required to truncate their output to {LINE\_MAX}. It is the responsibility of the

8416 application to deal with this. If the output of one of those utilities is to be piped into another

8417 of the standard utilities, line length restrictions will have to be considered; the *fold* utility,

8418 among others, could be used to ensure that only reasonable line lengths reach utilities or

8419 applications.

8420 {LINK\_MAX}

8421 This is defined by the System Interfaces volume of IEEE Std. 1003.1-200x.

8422 {MAX\_CANON}

8423 {MAX\_INPUT}

8424 {NAME\_MAX}

8425 {NGROUPS\_MAX}

8426 {OPEN\_MAX}

8427 {PATH\_MAX}

8428 {PIPE\_BUF}

8429 These limits are defined by the System Interfaces volume of IEEE Std. 1003.1-200x. Note that

8430 the byte lengths described by some of these values continue to represent bytes, even if the

8431 applicable character set uses a multi-byte encoding.

8432 {RE\_DUP\_MAX}  
 8433 The value selected is consistent with historical practice. Although the name implies that it  
 8434 applies to all REs, only BREs use the interval notation  $\{m,n\}$  addressed by this limit.

8435 {POSIX2\_SYMLINKS}  
 8436 The {POSIX2\_SYMLINKS} variable indicates that the underlying operating system supports  
 8437 the creation of symbolic links in specific directories. Many of the utilities defined in  
 8438 IEEE Std. 1003.1-200x that deal with symbolic links do not depend on this value. For  
 8439 example, a utility that follows symbolic links (or does not, as the case may be) will only be  
 8440 affected by a symbolic link if it encounters one. Presumably, a file system that does not  
 8441 support symbolic links will not contain any. This variable does affect such utilities as *ln -s*  
 8442 and *pax* that attempt to create symbolic links.

8443 {POSIX2\_SYMLINKS} was developed even though there is no comparable configuration  
 8444 value in the IEEE P1003.1a draft standard.

8445 There are different limits associated with command lines and input to utilities, depending on the  
 8446 method of invocation. In the case of a C program *exec*-ing a utility, {ARG\_MAX} is the  
 8447 underlying limit. In the case of the shell reading a script and *exec*-ing a utility, {LINE\_MAX}  
 8448 limits the length of lines the shell is required to process, and {ARG\_MAX} will still be a limit. If a  
 8449 user is entering a command on a terminal to the shell, requesting that it invoke the utility,  
 8450 {MAX\_INPUT} may restrict the length of the line that can be given to the shell to a value below  
 8451 {LINE\_MAX}.

8452 When an option is supported, *getconf* returns a value of 1. For example, when C development is  
 8453 supported:

```
8454     if [ "$(getconf POSIX2_C_DEV)" -eq 1 ]; then
8455         echo C supported
8456     fi
```

8457 The *sysconf()* function in the C-Language Binding option would return 1.

8458 The following comments describe selection criteria for the symbols and their values:

8459 POSIX2\_C\_BIND  
 8460 POSIX2\_C\_DEV  
 8461 POSIX2\_FORT\_DEV  
 8462 POSIX2\_FORT\_RUN  
 8463 POSIX2\_SW\_DEV  
 8464 POSIX2\_UPE

8465 It is possible for some (usually privileged) operations to remove utilities that support these  
 8466 options or otherwise to render these options unsupported. The header files, the *sysconf()*  
 8467 function, or the *getconf* utility will not necessarily detect such actions, in which case they  
 8468 should not be considered as rendering the implementation non-conforming. A test suite  
 8469 should not attempt tests such as:

```
8470     rm /usr/bin/c89
8471     getconf POSIX2_C_DEV
```

8472 POSIX2\_LOCALEDEF  
 8473 This symbol was introduced to allow implementations to restrict supported locales to only  
 8474 those supplied by the implementation.



8475 **C.1.10 Grammar Conventions**

8476 There is no additional rationale for this section.

8477 **C.1.11 Utility Description Defaults**8478 This section is arranged with headings in the same order as all the utility descriptions. It is a  
8479 collection of related and unrelated information concerning

8480 1. The default actions of utilities

8481 2. The meanings of notations used in IEEE Std. 1003.1-200x that are specific to individual  
8482 utility sections8483 Although this material may seem out of place here, it is important that this information appear  
8484 before any of the utilities to be described later.8485 **NAME**

8486 There is no additional rationale provided for this section.

8487 **SYNOPSIS**

8488 There is no additional rationale provided for this section.

8489 **DESCRIPTION**

8490 There is no additional rationale provided for this section.

8491 **OPTIONS**8492 Although it has not always been possible, the standard developers tried to avoid repeating  
8493 information to reduce the risk that duplicate explanations could each be modified differently.8494 The need to recognize `--` is required because portable applications need to shield their operands  
8495 from any arbitrary options that the implementation may provide as an extension. For example, if  
8496 the standard utility *foo* is listed as taking no options, and the application needed to give it a path  
8497 name with a leading hyphen, it could safely do it as:8498 `foo -- -myfile`8499 and avoid any problems with `-m` used as an extension.8500 **OPERANDS**8501 The usage of `-` is never shown in the SYNOPSIS. Similarly, the usage of `--` is never shown.8502 The requirement for processing operands in command-line order is to avoid a “WeirdNIX”  
8503 utility that might choose to sort the input files alphabetically, by size, or by directory order.  
8504 Although this might be acceptable for some utilities, in general the programmer has a right to  
8505 know exactly what order will be chosen.8506 Some of the standard utilities take multiple *file* operands and act as if they were processing the  
8507 concatenation of those files. For example:8508 `asa file1 file2`

8509 and:

8510 `cat file1 file2 | asa`

8511 have similar results when questions of file access, errors, and performance are ignored. Other  
8512 utilities such as *grep* or *wc* have completely different results in these two cases. This latter type of  
8513 utility is always identified in its DESCRIPTION or OPERANDS sections, whereas the former is  
8514 not. Although it might be possible to create a general assertion about the former case, the  
8515 following points must be addressed:

- 8516 • Access times for the files might be different in the operand case *versus* the *cat* case.
- 8517 • The utility may have error messages that are cognizant of the input file name, and this added  
8518 value should not be suppressed. (As an example, *awk* sets a variable with the file name at  
8519 each file boundary.)

## 8520 STDIN

8521 There is no additional rationale provided for this section.

## 8522 INPUT FILES

8523 A conforming application cannot assume the following three commands are equivalent:

```
8524 tail -n +2 file  
8525 (sed -n 1q; cat) < file  
8526 cat file | (sed -n 1q; cat)
```

8527 The second command is equivalent to the first only when the file is seekable. In the third  
8528 command, if the file offset in the open file description were not unspecified, *sed* would have to be  
8529 implemented so that it read from the pipe 1 byte at a time or it would have to employ some  
8530 method to seek backwards on the pipe. Such functionality is not defined currently in POSIX.1  
8531 and does not exist on all historical systems. Other utilities, such as *head*, *read*, and *sh*, have similar  
8532 properties, so the restriction is described globally in this section.

8533 The definition of *text file* is strictly enforced for input to the standard utilities; very few of them  
8534 list exceptions to the undefined results called for here. (Of course, “undefined” here does not  
8535 mean that historical implementations necessarily have to change to start indicating error  
8536 conditions. Conforming applications cannot rely on implementations succeeding or failing when  
8537 non-text files are used.)

8538 The utilities that allow line continuation are generally those that accept input languages, rather  
8539 than pure data. It would be unusual for an input line of this type to exceed {LINE\_MAX} bytes  
8540 and unreasonable to require that the implementation allow unlimited accumulation of multiple  
8541 lines, each of which could reach {LINE\_MAX}. Thus, for a portable application the total of all  
8542 the continued lines in a set cannot exceed {LINE\_MAX}.

8543 The format description is intended to be sufficiently rigorous to allow other applications to  
8544 generate these input files. However, since <blank>s can legitimately be included in some of the  
8545 fields described by the standard utilities, particularly in locales other than the POSIX locale, this  
8546 intent is not always realized.

## 8547 ENVIRONMENT VARIABLES

8548 There is no additional rationale provided for this section.

8549 **ASYNCHRONOUS EVENTS**

8550 Because there is no language prohibiting it, a utility is permitted to catch a signal, perform some  
8551 additional processing (such as deleting temporary files), restore the default signal action (or  
8552 action inherited from the parent process), and resignal itself.

8553 **STDOUT**

8554 The format description is intended to be sufficiently rigorous to allow post-processing of output  
8555 by other programs, particularly by an *awk* or *lex* parser.

8556 **STDERR**

8557 This section does not describe error messages that refer to incorrect operation of the utility.  
8558 Consider a utility that processes program source code as its input. This section is used to  
8559 describe messages produced by a correctly operating utility that encounters an error in the  
8560 program source code on which it is processing. However, a message indicating that the utility  
8561 had insufficient memory in which to operate would not be described.

8562 Some utilities have traditionally produced warning messages without returning a non-zero exit  
8563 status; these are specifically noted in their sections. Other utilities shall not write to standard  
8564 error if they complete successfully, unless the implementation provides some sort of extension  
8565 to increase the verbosity or debugging level.

8566 The format descriptions are intended to be sufficiently rigorous to allow post-processing of  
8567 output by other programs.

8568 **OUTPUT FILES**

8569 The format description is intended to be sufficiently rigorous to allow post-processing of output  
8570 by other programs, particularly by an *awk* or *lex* parser.

8571 Receipt of the SIGQUIT signal should generally cause termination (unless in some debugging  
8572 mode) that would bypass any attempted recovery actions.

8573 **EXTENDED DESCRIPTION**

8574 There is no additional rationale provided for this section.

8575 **EXIT STATUS**

8576 Note the additional discussion of exit values in *Exit Status for Commands* in the *sh* utility. It  
8577 describes requirements for returning exit values greater than 125.

8578 A utility may list zero as a successful return, 1 as a failure for a specific reason, and greater than  
8579 1 as “an error occurred”. In this case, unspecified conditions may cause a 2 or 3, or other value,  
8580 to be returned. A strictly conforming application should be written so that it tests for successful  
8581 exit status values (zero in this case), rather than relying upon the single specific error value listed  
8582 in IEEE Std. 1003.1-200x. In that way, it will have maximum portability, even on  
8583 implementations with extensions.

8584 The standard developers are aware that the general non-enumeration of errors makes it difficult  
8585 to write test suites that test the *incorrect* operation of utilities. There are some historical  
8586 implementations that have expended effort to provide detailed status messages and a helpful  
8587 environment to bypass or explain errors, such as prompting, retrying, or ignoring unimportant  
8588 syntax errors; other implementations have not. Since there is no realistic way to mandate system  
8589 behavior in cases of undefined application actions or system problems—in a manner acceptable  
8590 to all cultures and environments—attention has been limited to the correct operation of utilities

8591 by the conforming application. Furthermore, the portable application does not need detailed  
8592 information concerning errors that it caused through incorrect usage or that it cannot correct.

8593 There is no description of defaults for this section because all of the standard utilities specify  
8594 something (or explicitly state “Unspecified”) for exit status.

#### 8595 **CONSEQUENCES OF ERRORS**

8596 Several actions are possible when a utility encounters an error condition, depending on the  
8597 severity of the error and the state of the utility. Included in the possible actions of various  
8598 utilities are: deletion of temporary or intermediate work files; deletion of incomplete files; and  
8599 validity checking of the file system or directory.

8600 The text about recursive traversing is meant to ensure that utilities such as *find* process as many  
8601 files in the hierarchy as they can. They should not abandon all of the hierarchy at the first error  
8602 and resume with the next command-line operand, but should attempt to keep going.

#### 8603 **APPLICATION USAGE**

8604 This section provides additional caveats, issues, and recommendations to the developer.

#### 8605 **EXAMPLES**

8606 This section provides sample usage.

#### 8607 **RATIONALE**

8608 There is no additional rationale provided for this section.

#### 8609 **FUTURE DIRECTIONS**

8610 FUTURE DIRECTIONS sections act as pointers to related work that may impact the interface in  
8611 the future, and often cautions the developer to architect the code to account for a change in this  
8612 area. Note that a future directions statement should not be taken as a commitment to adopt a  
8613 feature or interface in the future.

#### 8614 **SEE ALSO**

8615 There is no additional rationale provided for this section.

#### 8616 **CHANGE HISTORY**

8617 There is no additional rationale provided for this section.

### 8618 **C.1.12 Considerations for Utilities in Support of Files of Arbitrary Size**

8619 This section is intended to clarify the requirements for utilities in support of large files.

8620 The utilities listed in this section are utilities which are used to perform administrative tasks  
8621 such as to create, move, copy, remove, change the permissions, or measure the resources of a  
8622 file. They are useful both as end-user tools and as utilities invoked by applications during  
8623 software installation and operation.

8624 The *chgrp*, *chmod*, *chown*, *ln*, and *rm* utilities probably require use of large file capable versions of  
8625 *stat()*, *lstat()*, *ftw()*, and the **stat** structure.

8626 The *cat*, *cksum*, *cmp*, *cp*, *dd*, *mv*, *sum*, and *touch* utilities probably require use of large file capable  
8627 versions of *creat()*, *open()*, and *fopen()*.

8628 The *cat*, *cksum*, *cmp*, *dd*, *df*, *du*, *ls*, and *sum* utilities may require writing large integer values. For  
8629 example:

- 8630 • The *cat* utility might have a `-n` option which counts <newline>s.
- 8631 • The *cksum* and *ls* utilities report file sizes.
- 8632 • The *cmp* utility reports the line number at which the first difference occurs, and also has a `-l`  
8633 option which reports file offsets.
- 8634 • The *dd*, *df*, *du*, *ls*, and *sum* utilities report block counts.

8635 The *dd*, *find*, and *test* utilities may need to interpret command arguments that contain 64-bit  
8636 values. For *dd*, the arguments include `skip=n`, `seek=n`, and `count=n`. For *find*, the arguments  
8637 include `-size n`. For *test*, the arguments are those associated with algebraic comparisons.

8638 The *df* utility might need to access large file systems with *statvfs()*.

8639 The *ulimit* utility will need to use large file capable versions of *getrlimit()* and *setrlimit()* and be  
8640 able to read and write large integer values.

8641 **C.2 Shell Command Language**8642 **C.2.1 Shell Introduction**

8643 The System V shell was selected as the starting point for the Shell and Utilities volume of  
 8644 IEEE Std. 1003.1-200x. The BSD C shell was excluded from consideration for the following  
 8645 reasons:

- 8646 • Most historically portable shell scripts assume the Version 7 Bourne shell, from which the  
 8647 System V shell is derived.
- 8648 • The majority of tutorial materials on shell programming assume the System V shell.

8649 The construct "#!" is reserved for implementations wishing to provide that extension. If it were  
 8650 not reserved, the Shell and Utilities volume of IEEE Std. 1003.1-200x would disallow it by forcing  
 8651 it to be a comment. As it stands, a POSIX-conforming application must not use "#!" as the first  
 8652 two characters of the file. An XSI-conforming application can use the construct "#!", since on  
 8653 XSI-conformant systems this is defined to denote an executable script, which matches historical  
 8654 practice. Invention of new meanings or extensions to the "#!" construct were rejected since they  
 8655 are beyond the scope of IEEE Std. 1003.1-200x.

8656 **C.2.2 Quoting**

8657 There is no additional rationale for this section.

8658 *C.2.2.1 Escape Character (Backslash)*

8659 There is no additional rationale for this section.

8660 *C.2.2.2 Single-Quotes*

8661 A backslash cannot be used to escape a single-quote in a single-quoted string. An embedded  
 8662 quote can be created by writing, for example: "'a'\ 'b'", which yields "a'b". (See the Shell  
 8663 and Utilities volume of IEEE Std. 1003.1-200x, Section 2.6.5, Field Splitting for a better  
 8664 understanding of how portions of words are either split into fields or remain concatenated.) A  
 8665 single token can be made up of concatenated partial strings containing all three kinds of quoting  
 8666 or escaping, thus permitting any combination of characters.

8667 *C.2.2.3 Double-Quotes*

8668 The escaped <newline> used for line continuation is removed entirely from the input and is not  
 8669 replaced by any white space. Therefore, it cannot serve as a token separator.

8670 In double-quoting, if a backslash is immediately followed by a character that would be  
 8671 interpreted as having a special meaning, the backslash is deleted and the subsequent character is  
 8672 taken literally. If a backslash does not precede a character that would have a special meaning, it  
 8673 is left in place unmodified and the character immediately following it is also left unmodified.  
 8674 Thus, for example:

8675       "\\$" → \$

8676       "a" → \a

8677 It would be desirable to include the statement “The characters from an enclosed "\${" to the  
 8678 matching '}' shall not be affected by the double quotes”, similar to the one for "\$()".  
 8679 However, historical practice in the System V shell prevents this.

8680 The requirement that double-quotes be matched inside "\${...}" within double-quotes and the  
 8681 rule for finding the matching '}' in the Shell and Utilities volume of IEEE Std. 1003.1-200x,  
 8682 Section 2.6.2, Parameter Expansion eliminate several subtle inconsistencies in expansion for  
 8683 historical shells in rare cases; for example:

```
8684     "${foo-bar}"
```

8685 yields **bar** when **foo** is not defined, and is an invalid substitution when **foo** is defined, in many  
 8686 historical shells. The differences in processing the "\${...}" form have led to inconsistencies  
 8687 between historical systems. A consequence of this rule is that single-quotes cannot be used to  
 8688 quote the '}' within "\${...}"; for example:

```
8689     unset bar  

  8690     foo="${bar-'}'"
```

8691 is invalid because the "\${...}" substitution contains an unpaired unescaped single-quote. The  
 8692 backslash can be used to escape the '}' in this example to achieve the desired result:

```
8693     unset bar  

  8694     foo="${bar-\\}'"
```

8695 The differences in processing the "\${...}" form have led to inconsistencies between the  
 8696 historical System V shell, BSD, and KornShells, and the text in the Shell and Utilities volume of  
 8697 IEEE Std. 1003.1-200x is an attempt to converge them without breaking too many applications.  
 8698 The only alternative to this compromise between shells would be to make the behavior  
 8699 unspecified whenever the literal characters ''', '{', '}', and '"' appear within "\${...}".  
 8700 To write a portable script that uses these values, a user would have to assign variables; for  
 8701 example:

```
8702     squote=\` dquote=\" lbrace='{` rbrace=}'`  

  8703     ${foo-`$squote$rbrace$squote`}
```

8704 rather than:

```
8705     ${foo-"' }' }
```

8706 Some systems have allowed the end of the word to terminate the backquoted command  
 8707 substitution, such as in:

```
8708     "`echo hello"
```

8709 This usage is undefined; the matching backquote is required by the Shell and Utilities volume of  
 8710 IEEE Std. 1003.1-200x. The other undefined usage can be illustrated by the example:

```
8711     sh -c '` echo "foo`'
```

8712 The description of the recursive actions involving command substitution can be illustrated with  
 8713 an example. Upon recognizing the introduction of command substitution, the shell parses input  
 8714 (in a new context), gathering the source for the command substitution until an unbalanced '}'  
 8715 or '`' is located. For example, in the following:

```
8716     echo "${date; echo "  

  8717         one" }"
```

8718 the double-quote following the *echo* does not terminate the first double-quote; it is part of the  
 8719 command substitution script. Similarly, in:

```
8720     echo "${echo *}"
```

8721 the asterisk is not quoted since it is inside command substitution; however:

8722           echo "\$ (echo "\*" )"

8723           is quoted (and represents the asterisk character itself).

### 8724 C.2.3 Token Recognition

8725           The "(" and ")" symbols are control operators in the KornShell, used for an alternative  
8726           syntax of an arithmetic expression command. A portable application cannot use "(" as a single  
8727           token (with the exception of the "\$ ( ( " form for shell arithmetic).

8728           The (3) rule about combining characters to form operators is not meant to preclude systems from  
8729           extending the shell language when characters are combined in otherwise invalid ways. Portable  
8730           applications cannot use invalid combinations, and test suites should not penalize systems that  
8731           take advantage of this fact. For example, the unquoted combination "|&" is not valid in a POSIX  
8732           script, but has a specific KornShell meaning.

8733           The (10) rule about '#' as the current character is the first in the sequence in which a new token  
8734           is being assembled. The '#' starts a comment only when it is at the beginning of a token. This  
8735           rule is also written to indicate that the search for the end-of-comment does not consider escaped  
8736           <newline> specially, so that a comment cannot be continued to the next line.

#### 8737 C.2.3.1 Alias Substitution

8738           The alias capability was added in the UPE because it is widely used in historical  
8739           implementations by interactive users.

8740           The definition of *alias name* precludes an alias name containing a slash character. Since the text  
8741           applies to the command words of simple commands, reserved words (in their proper places)  
8742           cannot be confused with aliases.

8743           The placement of alias substitution in token recognition makes it clear that it precedes all of the  
8744           word expansion steps.

8745           An example concerning trailing <blank> characters and reserved words follows. If the user  
8746           types:

8747           \$ alias foo="/bin/ls "

8748           \$ alias while="/"

8749           The effect of executing:

8750           \$ while true

8751           > do

8752           > echo "Hello, World"

8753           > done

8754           is a never-ending sequence of "Hello, World" strings to the screen. However, if the user  
8755           types:

8756           \$ foo while

8757           the result is an *ls* listing of /. Since the alias substitution for **foo** ends in a <space> character, the  
8758           next word is checked for alias substitution. The next word, **while**, has also been aliased, so it is  
8759           substituted as well. Since it is not in the proper position as a command word, it is not recognized  
8760           as a reserved word.

8761           If the user types:

8762           \$ foo; while



8763 **while** retains its normal reserved-word properties.

#### 8764 **C.2.4 Reserved Words**

8765 All reserved words are recognized syntactically as such in the contexts described. However, note  
8766 that **in** is the only meaningful reserved word after a **case** or **for**; similarly, **in** is not meaningful as  
8767 the first word of a simple command.

8768 Reserved words are recognized only when they are delimited (that is, meet the definition of the  
8769 Base Definitions volume of IEEE Std. 1003.1-200x, Section 3.437, Word), whereas operators are  
8770 themselves delimiters. For instance, '( ' and ' ) ' are control operators, so that no <space>  
8771 character is needed in (*list*). However, '{ ' and ' } ' are reserved words in {*list*}, so that in this  
8772 case the leading <space> character and semicolon are required.

8773 The list of unspecified reserved words is from the KornShell, so portable applications cannot use  
8774 them in places a reserved word would be recognized. This list contained **time** in early proposals,  
8775 but it was removed when the *time* utility was selected for the Shell and Utilities volume of  
8776 IEEE Std. 1003.1-200x.

8777 There was a strong argument for promoting braces to operators (instead of reserved words), so  
8778 they would be syntactically equivalent to subshell operators. Concerns about compatibility  
8779 outweighed the advantages of this approach. Nevertheless, portable applications should  
8780 consider quoting '{ ' and ' } ' when they represent themselves.

8781 The restriction on ending a name with a colon is to allow future implementations that support  
8782 named labels for flow control; see the RATIONALE for the *break* built-in utility .

8783 It is possible that a future version of the Shell and Utilities volume of IEEE Std. 1003.1-200x may  
8784 require that '{ ' and ' } ' be treated individually as control operators, although the token "{ }"  
8785 will probably be a special-case exemption from this because of the often-used *find*{ } construct.

#### 8786 **C.2.5 Parameters and Variables**

##### 8787 *C.2.5.1 Positional Parameters*

8788 There is no additional rationale for this section.

##### 8789 *C.2.5.2 Special Parameters*

8790 Most historical implementations implement subshells by forking; thus, the special parameter  
8791 '\$ ' does not necessarily represent the process ID of the shell process executing the commands  
8792 since the subshell execution environment preserves the value of '\$ '.

8793 If a subshell were to execute a background command, the value of "\$ ! " for the parent would  
8794 not change. For example:

```
8795     (
8796     date &
8797     echo $!
8798     )
8799     echo $!
```

8800 would echo two different values for "\$ ! " .

8801 The "\$ - " special parameter can be used to save and restore *set* options:

```

8802     Save=$(echo $- | sed 's/[ics]//g')
8803     ...
8804     set +aCefnuvx
8805     if [ -n "$Save" ]; then
8806         set -$Save
8807     fi

```

8808 The three options are removed using *sed* in the example because they may appear in the value of  
8809 "\$-" (from the *sh* command line), but are not valid options to *set*.

8810 The descriptions of parameters '\*' and '@' assume the reader is familiar with the field splitting  
8811 discussion in the Shell and Utilities volume of IEEE Std. 1003.1-200x, Section 2.6.5, Field Splitting  
8812 and understands that portions of the word remain concatenated unless there is some reason to  
8813 split them into separate fields.

8814 Some examples of the '\*' and '@' properties, including the concatenation aspects:

```

8815     set "abc" "def ghi" "jkl"
8816     echo $*      => "abc" "def" "ghi" "jkl"
8817     echo "$*"    => "abc def ghi jkl"
8818     echo $@      => "abc" "def" "ghi" "jkl"

```

8819 but:

```

8820     echo "$@"    => "abc" "def ghi" "jkl"
8821     echo "xx$@yy" => "xxabc" "def ghi" "jklyy"
8822     echo "$@$@"  => "abc" "def ghi" "jklabc" "def ghi" "jkl"

```

8823 In the preceding examples, the double-quote characters that appear after the "=>" do not appear  
8824 in the output and are used only to illustrate word boundaries.

8825 The following example illustrates the effect of setting *IFS* to a null string:

```

8826     $ IFS=''
8827     $ set foo bar bam
8828     $ echo "$@"
8829     foo bar bam
8830     $ echo "$*"
8831     foobarbam
8832     $ unset IFS
8833     $ echo "$*"
8834     foo bar bam

```

### 8835 C.2.5.3 Shell Variables

8836 See the discussion of *IFS* in Section C.2.6.5 (on page 3532).

8837 The prohibition on *LC\_CTYPE* changes affecting lexical processing protects the shell  
8838 implementor (and the shell programmer) from the ill effects of changing the definition of  
8839 <blank> or the set of alphabetic characters in the current environment. It would probably not be  
8840 feasible to write a compiled version of a shell script without this rule. The rule applies only to  
8841 the current invocation of the shell and its subshells—invoking a shell script or performing *exec sh*  
8842 would subject the new shell to the changes in *LC\_CTYPE*.

8843 Other common environment variables used by historical shells are not specified by the Shell and  
8844 Utilities volume of IEEE Std. 1003.1-200x, but they should be reserved for the historical uses.

8845		Tilde expansion for components of the <i>PATH</i> in an assignment such as:
8846		<code>PATH=~hlj/bin:~dwc/bin:\$PATH</code>
8847		is a feature of some historical shells and is allowed by the wording of the Shell and Utilities
8848		volume of IEEE Std. 1003.1-200x, Section 2.6.1, Tilde Expansion. Note that the tildes are
8849		expanded during the assignment to <i>PATH</i> , not when <i>PATH</i> is accessed during command search.
8850		The following entries represent additional information about variables included in the Shell and
8851		Utilities volume of IEEE Std. 1003.1-200x, or rationale for common variables in use by shells that
8852		have been excluded:
8853	—	(Underscore.) While underscore is historical practice, its overloaded usage in
8854		the KornShell is confusing, and it has been omitted from the Shell and Utilities
8855		volume of IEEE Std. 1003.1-200x.
8856	<i>ENV</i>	This variable can be used to set aliases and other items local to the invocation
8857		of a shell. The file referred to by <i>ENV</i> differs from <i>\$HOME/.profile</i> in that
8858		<i>.profile</i> is typically executed at session start-up, whereas the <i>ENV</i> file is
8859		executed at the beginning of each shell invocation. The <i>ENV</i> value is
8860		interpreted in a manner similar to a dot script, in that the commands are
8861		executed in the current environment and the file needs to be readable, but not
8862		executable. However, unlike dot scripts, no <i>PATH</i> searching is performed.
8863		This is used as a guard against Trojan Horse security breaches.
8864	<i>ERRNO</i>	This variable was omitted from the Shell and Utilities volume of
8865		IEEE Std. 1003.1-200x because the values of error numbers are not defined in
8866		IEEE Std. 1003.1-200x in a portable manner.
8867	<i>FCEDIT</i>	Since this variable affects only the <i>fc</i> utility, it has been omitted from this more
8868		global place. The value of <i>FCEDIT</i> does not affect the command line editing
8869		mode in the shell; see the description of <i>set -o vi</i> in the <i>set</i> built-in utility.
8870	<i>PS1</i>	This variable is used for interactive prompts. Historically, the “superuser”
8871		has had a prompt of ‘#’. Since privileges are not required to be monolithic, it
8872		is difficult to define which privileges should cause the alternate prompt.
8873		However, a sufficiently powerful user should be reminded of that power by
8874		having an alternate prompt.
8875	<i>PS3</i>	This variable is used by the KornShell for the <i>select</i> command. Since the POSIX
8876		shell does not include <i>select</i> , <i>PS3</i> was omitted.
8877	<i>PS4</i>	This variable is used for shell debugging. For example, the following script:
8878		<code>PS4=' [ \${LINENO} ]+ '</code>
8879		<code>set -x</code>
8880		<code>echo Hello</code>
8881		writes the following to standard error:
8882		<code>[3]+ echo Hello</code>
8883	<i>RANDOM</i>	This pseudo-random number generator was not seen as being useful to
8884		interactive users.
8885	<i>SECONDS</i>	Although this variable is sometimes used with <i>PS1</i> to allow the display of the
8886		current time in the prompt of the user, it is not one that would be manipulated
8887		frequently enough by an interactive user to include in the Shell and Utilities
8888		volume of IEEE Std. 1003.1-200x.

8889 **C.2.6 Word Expansions**

8890 Step (2) refers to the “portions of fields generated by step (1)”. For example, if the word being  
 8891 expanded were "\$x+\$y" and *IFS*=+, the word would be split only if "\$x" or "\$y" contained  
 8892 '+'; the '+' in the original word was not generated by step (1).

8893 *IFS* is used for performing field splitting on the results of parameter and command substitution;  
 8894 it is not used for splitting all fields. Previous versions of the shell used it for splitting all fields  
 8895 during field splitting, but this has severe problems because the shell can no longer parse its own  
 8896 script. There are also important security implications caused by this behavior. All useful  
 8897 applications of *IFS* use it for parsing input of the *read* utility and for splitting the results of  
 8898 parameter and command substitution.

8899 The rule concerning expansion to a single field requires that if **foo=abc** and **bar=def**, that:

8900 "\$foo" "\$bar"

8901 expands to the single field:

8902 abcdef

8903 The rule concerning empty fields can be illustrated by:

```
8904 $ unset foo
8905 $ set $foo bar ' ' xyz "$foo" abc
8906 $ for i
8907 > do
8908 >     echo "-$i-"
8909 > done
8910 -bar-
8911 --
8912 -xyz-
8913 --
8914 -abc-
```

8915 Step (1) indicates that parameter expansion, command substitution, and arithmetic expansion  
 8916 are all processed simultaneously as they are scanned. For example, the following is valid  
 8917 arithmetic:

```
8918 x=1
8919 echo $(( $(echo 3)+$x ))
```

8920 An early proposal stated that tilde expansion preceded the other steps, but this is not the case in  
 8921 known historical implementations; if it were, and if a referenced home directory contained a '\$'  
 8922 character, expansions would result within the directory name.

8923 **C.2.6.1 Tilde Expansion**

8924 Tilde expansion generally occurs only at the beginning of words, but an exception based on  
 8925 historical practice has been included:

```
8926 PATH=/posix/bin:~djk/bin
```

8927 This is eligible for tilde expansion because tilde follows a colon and none of the relevant  
 8928 characters is quoted. Consideration was given to prohibiting this behavior because any of the  
 8929 following are reasonable substitutes:

```
8930 PATH=$(printf %s ~karels/bin : ~bostic/bin)
```

```

8931     for Dir in ~maat/bin ~srb/bin ...
8932     do
8933         PATH=${PATH:+$PATH:}$Dir
8934     done

```

8935 In the first command, explicit colons are used for each directory. In all cases, the shell performs  
 8936 tilde expansion on each directory because all are separate words to the shell.

8937 Note that expressions in operands such as:

```

8938     make -k mumble LIBDIR=~chet/lib

```

8939 do not qualify as shell variable assignments, and tilde expansion is not performed (unless the  
 8940 command does so itself, which *make* does not).

8941 Because of the requirement that the word is not quoted, the following are not equivalent; only  
 8942 the last causes tilde expansion:

```

8943     \~hlj/   ~h\lj/   ~"hlj"/   ~hlj\ /   ~hlj/

```

8944 In an early proposal, tilde expansion occurred following any unquoted equals sign or colon, but  
 8945 this was removed because of its complexity and to avoid breaking commands such as:

```

8946     rcp hostname:~marc/.profile .

```

8947 A suggestion was made that the special sequence "\$~" should be allowed to force tilde  
 8948 expansion anywhere. Since this is not historical practice, it has been left for future  
 8949 implementations to evaluate. (The description in the Shell and Utilities volume of  
 8950 IEEE Std. 1003.1-200x, Section 2.2, Quoting requires that a dollar sign be quoted to represent  
 8951 itself, so the "\$~" combination is already unspecified.)

8952 The results of giving tilde with an unknown login name are undefined because the KornShell  
 8953 "~+" and "~-" constructs make use of this condition, but in general it is an error to give an  
 8954 incorrect login name with tilde. The results of having *HOME* unset are unspecified because some  
 8955 historical shells treat this as an error.

#### 8956 C.2.6.2 Parameter Expansion

8957 The rule for finding the closing '}' in "\${...}" is the one used in the KornShell and is  
 8958 upwardly-compatible with the Bourne shell, which does not determine the closing '}' until the  
 8959 word is expanded. The advantage of this is that incomplete expansions, such as:

```

8960     ${foo

```

8961 can be determined during tokenization, rather than during expansion.

8962 The string length and substring capabilities were included because of the demonstrated need for  
 8963 them, based on their usage in other shells, such as C shell and KornShell.

8964 Historical versions of the KornShell have not performed tilde expansion on the word part of  
 8965 parameter expansion; however, it is more consistent to do so.

#### 8966 C.2.6.3 Command Substitution

8967 The "\$()" form of command substitution solves a problem of inconsistent behavior when using  
 8968 backquotes. For example:

8969  
8970  
8971  
8972  
8973

Command	Output
echo '\\$x'	\\$x
echo `echo '\\$x'`	\$x
echo \$(echo '\\$x')	\\$x

8974  
8975  
8976  
8977

Additionally, the backquoted syntax has historical restrictions on the contents of the embedded command. While the newer "\$()" form can process any kind of valid embedded script, the backquoted form cannot handle some valid scripts that include backquotes. For example, these otherwise valid embedded scripts do not work in the left column, but do work on the right:

8978  
8979  
8980  
8981  
8982  
  
8983  
8984  
8985  
  
8986  
8987  
8988

```

echo `
cat <<\eof
a here-doc with `
eof
`

echo `
echo abc # a comment with `
`

echo `
echo ` `
`

echo $(
cat <<\eof
a here-doc with )
eof
)

echo $(
echo abc # a comment with )
)

echo $(
echo ` `)
)

```

8989  
8990  
8991

Because of these inconsistent behaviors, the backquoted variety of command substitution is not recommended for new applications that nest command substitutions or attempt to embed complex scripts.

8992

The KornShell feature:

8993  
8994

If *command* is of the form `<word`, *word* is expanded to generate a path name, and the value of the command substitution is the contents of this file with any trailing `<newline>`s deleted.

8995  
8996  
8997

was omitted from the Shell and Utilities volume of IEEE Std. 1003.1-200x because `$(cat word)` is an appropriate substitute. However, to prevent breaking numerous scripts relying on this feature, it is unspecified to have a script within `"$()"` that has only redirections.

8998  
8999

The requirement to separate `"$("` and `'('` when a single subshell is command-substituted is to avoid any ambiguities with arithmetic expansion.

#### 9000 C.2.6.4 Arithmetic Expansion

9001  
9002  
9003  
9004  
9005

The `"(( ))"` form of KornShell arithmetic in early proposals was omitted. The standard developers concluded that there was a strong desire for some kind of arithmetic evaluator to replace *expr*, and that relating it to `'$'` makes it work well with the standard shell language, and it provides access to arithmetic evaluation in places where accessing a utility would be inconvenient.

9006  
9007  
9008  
9009  
9010  
9011  
9012  
9013  
9014

The syntax and semantics for arithmetic were changed for the ISO/IEC 9945-2:1993 standard. The language is essentially a pure arithmetic evaluator of constants and operators (excluding assignment) and represents a simple subset of the previous arithmetic language (which was derived from the KornShell `"(( ))"` construct). The syntax was changed from that of a command denoted by `((expression))` to an expansion denoted by `$((expression))`. The new form is a dollar expansion `'$'` that evaluates the expression and substitutes the resulting value. Objections to the previous style of arithmetic included that it was too complicated, did not fit in well with the use of variables in the shell, and its syntax conflicted with subshells. The justification for the new syntax is that the shell is traditionally a macro language, and if a new

9015 feature is to be added, it should be accomplished by extending the capabilities presented by the  
 9016 current model of the shell, rather than by inventing a new one outside the model; adding a new  
 9017 dollar expansion was perceived to be the most intuitive and least destructive way to add such a  
 9018 new capability.

9019 In early proposals, a form  $\$(expression)$  was used. It was functionally equivalent to the " $\$(())$ "  
 9020 of the current text, but objections were lodged that the 1988 KornShell had already implemented  
 9021 " $\$(())$ " and there was no compelling reason to invent yet another syntax. Furthermore, the  
 9022 " $\$[]$ " syntax had a minor incompatibility involving the patterns in **case** statements.

9023 The portion of the ISO C standard arithmetic operations selected corresponds to the operations  
 9024 historically supported in the KornShell.

9025 It was concluded that the *test* command (**D**) was sufficient for the majority of relational arithmetic  
 9026 tests, and that tests involving complicated relational expressions within the shell are rare, yet  
 9027 could still be accommodated by testing the value of " $\$(())$ " itself. For example:

```
9028     # a complicated relational expression
9029     while [  $\$( (( $x + $y ) / ( $a * $b )) < ( $foo * $bar )) -ne 0 ]$ 
```

9030 or better yet, the rare script that has many complex relational expressions could define a  
 9031 function like this:

```
9032     val() {
9033         return  $\$( ! $1 )$ 
9034     }
```

9035 and complicated tests would be less intimidating:

```
9036     while val  $\$( (( $x + $y ) / ( $a * $b )) < ( $foo * $bar ))$ 
9037     do
9038         # some calculations
9039     done
```

9040 A suggestion that was not adopted was to modify *true* and *false* to take an optional argument,  
 9041 and *true* would exit true only if the argument was non-zero, and *false* would exit false only if the  
 9042 argument was non-zero:

```
9043     while true  $\$( ($x > 5 && $y <= 25) )$ 
```

9044 There is a minor portability concern with the new syntax. The example  $\$(2+2)$  could have been  
 9045 intended to mean a command substitution of a utility named *2+2* in a subshell. The standard  
 9046 developers considered this to be obscure and isolated to some KornShell scripts (because " $\$( )$ "  
 9047 command substitution existed previously only in the KornShell). The text on command  
 9048 substitution requires that the " $\$($ " and ' $($ ' be separate tokens if this usage is needed.

9049 An example such as:

```
9050     echo  $\$( (echo hi); (echo there) )$ 
```

9051 should not be misinterpreted by the shell as arithmetic because attempts to balance the  
 9052 parentheses pairs would indicate that they are subshells. However, as indicated by the Base  
 9053 Definitions volume of IEEE Std. 1003.1-200x, Section 3.115, Control Operator, a conforming  
 9054 application must separate two adjacent parentheses with white space to indicate nested  
 9055 subshells.

9056 **C.2.6.5 Field Splitting**

9057 The operation of field splitting using *IFS*, as described in early proposals, was based on the way  
 9058 the KornShell splits words, but it is incompatible with other common versions of the shell.  
 9059 However, each has merit, and so a decision was made to allow both. If the *IFS* variable is unset  
 9060 or is `<space><tab><newline>`, the operation is equivalent to the way the System V shell splits  
 9061 words. Using characters outside the `<space><tab><newline>` set yields the KornShell behavior,  
 9062 where each of the non-`<space><tab><newline>` characters is significant. This behavior, which  
 9063 affords the most flexibility, was taken from the way the original *awk* handled field splitting.

9064 Rule (3) can be summarized as a pseudo-ERE:

9065  $(s^*ns^*|s^+)$

9066 where *s* is an *IFS* white space character and *n* is a character in the *IFS* that is not white space.  
 9067 Any string matching that ERE delimits a field, except that the *s+* form does not delimit fields at  
 9068 the beginning or the end of a line. For example, if *IFS* is `<space>/<comma>/<tab>`, the string:

9069 `<space><space>red<space><space>, <space>white<space>blue`

9070 yields the three colors as the delimited fields.

9071 **C.2.6.6 Path Name Expansion**

9072 There is no additional rationale for this section.

9073 **C.2.6.7 Quote Removal**

9074 There is no additional rationale for this section.

9075 **C.2.7 Redirection**

9076 In the System Interfaces volume of IEEE Std. 1003.1-200x, file descriptors are integers in the  
 9077 range 0–(`{OPEN_MAX}`–1). The file descriptors discussed in the Shell and Utilities volume of  
 9078 IEEE Std. 1003.1-200x, Section 2.7, Redirection are that same set of small integers.

9079 Having multi-digit file descriptor numbers for I/O redirection can cause some obscure  
 9080 compatibility problems. Specifically, scripts that depend on an example command:

9081 `echo 22>/dev/null`

9082 echoing 2 to standard error or 22 to standard output are no longer portable. However, the file  
 9083 descriptor number still must be delimited from the preceding text. For example:

9084 `cat file2>foo`

9085 writes the contents of **file2**, not the contents of **file**.

9086 The "`>`" format of output redirection was adopted from the KornShell. Along with the  
 9087 *noclobber* option, `set -C`, it provides a safety feature to prevent inadvertent overwriting of  
 9088 existing files. (See the RATIONALE for the *pathchk* utility for why this step was taken.) The  
 9089 restriction on regular files is historical practice.

9090 The System V shell and the KornShell have differed historically on path name expansion of *word*;  
 9091 the former never performed it, the latter only when the result was a single field (file). As a  
 9092 compromise, it was decided that the KornShell functionality was useful, but only as a shorthand  
 9093 device for interactive users. No reasonable shell script would be written with a command such  
 9094 as:

9095 `cat foo > a*`



9096 Thus, shell scripts are prohibited from doing it, while interactive users can select the shell with  
9097 which they are most comfortable.

9098 The construct `2>&1` is often used to redirect standard error to the same file as standard output.  
9099 Since the redirections take place beginning to end, the order of redirections is significant. For  
9100 example:

```
9101     ls > foo 2>&1
```

9102 directs both standard output and standard error to file **foo**. However:

```
9103     ls 2>&1 > foo
```

9104 only directs standard output to file **foo** because standard error was duplicated as standard  
9105 output before standard output was directed to file **foo**.

9106 The "<>" operator could be useful in writing an application that worked with several terminals,  
9107 and occasionally wanted to start up a shell. That shell would in turn be unable to run  
9108 applications that run from an ordinary controlling terminal unless it could make use of "<>"  
9109 redirection. The specific example is a historical version of the pager *more*, which reads from  
9110 standard error to get its commands, so standard input and standard output are both available  
9111 for their usual usage. There is no way of saying the following in the shell without "<>":

```
9112     cat food | more - >/dev/tty03 2<>/dev/tty03
```

9113 Another example of "<>" is one that opens `/dev/tty` on file descriptor 3 for reading and writing:

```
9114     exec 3<> /dev/tty
```

9115 An example of creating a lock file for a critical code region:

```
9116     set -C
9117     until    2> /dev/null > lockfile
9118     do      sleep 30
9119     done
9120     set +C
9121     perform critical function
9122     rm lockfile
```

9123 Since `/dev/null` is not a regular file, no error is generated by redirecting to it in *noclobber* mode.

9124 Tilde expansion is not performed on a here-document because the data is treated as if it were  
9125 enclosed in double quotes.

#### 9126 *C.2.7.1 Redirecting Input*

9127 There is no additional rationale for this section.

#### 9128 *C.2.7.2 Redirecting Output*

9129 There is no additional rationale for this section.

#### 9130 *C.2.7.3 Appending Redirected Output*

9131 There is no additional rationale for this section.

9132 C.2.7.4 *Here-Document*

9133 There is no additional rationale for this section.

9134 C.2.7.5 *Duplicating an Input File Descriptor*

9135 There is no additional rationale for this section.

9136 C.2.7.6 *Duplicating an Output File Descriptor*

9137 There is no additional rationale for this section.

9138 C.2.7.7 *Open File Descriptors for Reading and Writing*

9139 There is no additional rationale for this section.

9140 **C.2.8 Exit Status and Errors**9141 C.2.8.1 *Consequences of Shell Errors*

9142 There is no additional rationale for this section.

9143 C.2.8.2 *Exit Status for Commands*

9144 There is a historical difference in *sh* and *ksh* non-interactive error behavior. When a command  
9145 named in a script is not found, some implementations of *sh* exit immediately, but *ksh* continues  
9146 with the next command. Thus, the Shell and Utilities volume of IEEE Std. 1003.1-200x says that  
9147 the shell “may” exit in this case. This puts a small burden on the programmer, who has to test  
9148 for successful completion following a command if it is important that the next command not be  
9149 executed if the previous command was not found. If it is important for the command to have  
9150 been found, it was probably also important for it to complete successfully. The test for successful  
9151 completion would not need to change.

9152 Historically, shells have returned an exit status of  $128+n$ , where  $n$  represents the signal number.  
9153 Since signal numbers are not standardized, there is no portable way to determine which signal  
9154 caused the termination. Also, it is possible for a command to exit with a status in the same range  
9155 of numbers that the shell would use to report that the command was terminated by a signal.  
9156 Implementations are encouraged to choose exit values greater than 256 to indicate programs  
9157 that terminate by a signal so that the exit status cannot be confused with an exit status generated  
9158 by a normal termination.

9159 Historical shells make the distinction between “utility not found” and “utility found but cannot  
9160 execute” in their error messages. By specifying two seldomly used exit status values for these  
9161 cases, 127 and 126 respectively, this gives an application the opportunity to make use of this  
9162 distinction without having to parse an error message that would probably change from locale to  
9163 locale. The *command*, *env*, *nohup*, and *xargs* utilities in the Shell and Utilities volume of  
9164 IEEE Std. 1003.1-200x have also been specified to use this convention.

9165 When a command fails during word expansion or redirection, most historical implementations  
9166 exit with a status of 1. However, there was some sentiment that this value should probably be  
9167 much higher so that an application could distinguish this case from the more normal exit status  
9168 values. Thus, the language “greater than zero” was selected to allow either method to be  
9169 implemented.

9170 **C.2.9 Shell Commands**

9171 A description of an “empty command” was removed from an early proposal because it is only  
 9172 relevant in the cases of *sh -c " "*, *system(" ")*, or an empty shell-script file (such as the  
 9173 implementation of *true* on some historical systems). Since it is no longer mentioned in the Shell  
 9174 and Utilities volume of IEEE Std. 1003.1-200x, it falls into the silently unspecified category of  
 9175 behavior where implementations can continue to operate as they have historically, but  
 9176 conforming applications do not construct empty commands. (However, note that *sh* does  
 9177 explicitly state an exit status for an empty string or file.) In an interactive session or a script with  
 9178 other commands, extra <newline>s or semicolons, such as;

```
9179     $ false
9180     $
9181     $ echo $?
9182     1
```

9183 would not qualify as the empty command described here because they would be consumed by  
 9184 other parts of the grammar.

9185 **C.2.9.1 Simple Commands**

9186 The enumerated list is used only when the command is actually going to be executed. For  
 9187 example, in:

```
9188     true || $foo *
```

9189 no expansions are performed.

9190 The following example illustrates both how a variable assignment without a command name  
 9191 affects the current execution environment, and how an assignment with a command name only  
 9192 affects the execution environment of the command:

```
9193     $ x=red
9194     $ echo $x
9195     red
9196     $ export x
9197     $ sh -c 'echo $x'
9198     red
9199     $ x=blue sh -c 'echo $x'
9200     blue
9201     $ echo $x
9202     red
```

9203 This next example illustrates that redirections without a command name are still performed:

```
9204     $ ls foo
9205     ls: foo: no such file or directory
9206     $ > foo
9207     $ ls foo
9208     foo
```

9209 A command without a command name, but one that includes a command substitution, has an  
 9210 exit status of the last command substitution that the shell performed. For example:

```
9211     if      x=$(command)
9212     then    ...
9213     fi
```

9214 An example of redirections without a command name being performed in a subshell shows that  
 9215 the here-document does not disrupt the standard input of the **while** loop:

```
9216     IFS=:
9217     while read a b
9218     do      echo $a
9219           <<-eof
9220           Hello
9221           eof
9222     done </etc/passwd
```

9223 Some examples of commands without command names in AND-OR lists:

```
9224     > foo || {
9225         echo "error: foo cannot be created" >&2
9226         exit 1
9227     }
9228     # set saved if /vmunix.save exists
9229     test -f /vmunix.save && saved=1
```

9230 Command substitution and redirections without command names both occur in subshells, but  
 9231 they are not necessarily the same ones. For example, in:

```
9232     exec 3> file
9233     var=$(echo foo >&3) 3>&1
```

9234 it is unspecified whether **foo** is echoed to the file or to standard output.

### 9235 **Command Search and Execution**

9236 This description requires that the shell can execute shell scripts directly, even if the underlying  
 9237 system does not support the common "#!" interpreter convention. That is, if file **foo** contains  
 9238 shell commands and is executable, the following executes **foo**:

```
9239     ./foo
```

9240 The command search shown here does not match all historical implementations. A more typical  
 9241 sequence has been:

- 9242 • Any built-in (special or regular)
- 9243 • Functions
- 9244 • Path search for executable files

9245 But there are problems with this sequence. Since the programmer has no idea in advance which  
 9246 utilities might have been built into the shell, a function cannot be used to override portably a  
 9247 utility of the same name. (For example, a function named *cd* cannot be written for many  
 9248 historical systems.) Furthermore, the *PATH* variable is partially ineffective in this case, and only  
 9249 a path name with a slash can be used to ensure a specific executable file is invoked.

9250 After the *execve()* failure described, the shell normally executes the file as a shell script. Some  
 9251 implementations, however, attempt to detect whether the file is actually a script and not an  
 9252 executable from some other architecture. The method used by the KornShell is allowed by the  
 9253 text that indicates non-text files may be bypassed.

9254 The sequence selected for the Shell and Utilities volume of IEEE Std. 1003.1-200x acknowledges  
 9255 that special built-ins cannot be overridden, but gives the programmer full control over which  
 9256 versions of other utilities are executed. It provides a means of suppressing function lookup (via

9257 the *command* utility) for the user's own functions and ensures that any regular built-ins or  
 9258 functions provided by the implementation are under the control of the path search. The  
 9259 mechanisms for associating built-ins or functions with executable files in the path are not  
 9260 specified by the Shell and Utilities volume of IEEE Std. 1003.1-200x, but the wording requires  
 9261 that if either is implemented, the application is not able to distinguish a function or built-in from  
 9262 an executable (other than in terms of performance, presumably). The implementation ensures  
 9263 that all effects specified by the Shell and Utilities volume of IEEE Std. 1003.1-200x resulting from  
 9264 the invocation of the regular built-in or function (interaction with the environment, variables,  
 9265 traps, and so on) are identical to those resulting from the invocation of an executable file.

### 9266 Examples

9267 Consider three versions of the *ls* utility:

- 9268 1. The application includes a shell function named *ls*.
- 9269 2. The user writes a utility named *ls* and puts it in **/fred/bin**.
- 9270 3. The example implementation provides *ls* as a regular shell built-in that is invoked (either  
 9271 by the shell or directly by *exec*) when the path search reaches the directory **/posix/bin**.

9272 If *PATH*=**/posix/bin**, various invocations yield different versions of *ls*:

9273

9274

9275

9276

9277

9278

9279

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

### 9280 C.2.9.2 Pipelines

9281 Because pipeline assignment of standard input or standard output or both takes place before  
 9282 redirection, it can be modified by redirection. For example:

```
9283 $ command1 2>&1 | command2
```

9284 sends both the standard output and standard error of *command1* to the standard input of  
 9285 *command2*.

9286 The reserved word **!** allows more flexible testing using AND and OR lists.

9287 It was suggested that it would be better to return a non-zero value if any command in the  
 9288 pipeline terminates with non-zero status (perhaps the bitwise-inclusive OR of all return values).  
 9289 However, the choice of the last-specified command semantics are historical practice and would  
 9290 cause applications to break if changed. An example of historical behavior:

```
9291 $ sleep 5 | (exit 4)
```

```
9292 $ echo $?
```

```
9293 4
```

```
9294 $ (exit 4) | sleep 5
```

```
9295 $ echo $?
```

```
9296 0
```

9297 *C.2.9.3 Lists*

9298 The equal precedence of "&&" and "||" is historical practice. The standard developers  
 9299 evaluated the model used more frequently in high-level programming languages, such as C, to  
 9300 allow the shell logical operators to be used for complex expressions in an unambiguous way, but  
 9301 they could not allow historical scripts to break in the subtle way unequal precedence might  
 9302 cause. Some arguments were posed concerning the "{" or "(" groupings that are required  
 9303 historically. There are some disadvantages to these groupings:

- 9304 • The "(" can be expensive, as they spawn other processes on some systems. This  
 9305 performance concern is primarily an implementation issue.
- 9306 • The "{" braces are not operators (they are reserved words) and require a trailing space  
 9307 after each '{', and a semicolon before each '}'. Most programmers (and certainly  
 9308 interactive users) have avoided braces as grouping constructs because of the problematic  
 9309 syntax required. Braces were not changed to operators because that would generate  
 9310 compatibility issues even greater than the precedence question; braces appear outside the  
 9311 context of a keyword in many shell scripts.

9312 **Asynchronous Lists**

9313 The grammar treats a construct such as:

```
9314     foo & bar & bam &
```

9315 as one "asynchronous list", but since the status of each element is tracked by the shell, the term  
 9316 "element of an asynchronous list" was introduced to identify just one of the **foo**, **bar**, or **bam**  
 9317 portions of the overall list.

9318 Unless the implementation has an internal limit, such as {CHILD\_MAX}, on the retained process  
 9319 IDs, it would require unbounded memory for the following example:

```
9320     while true
9321     do         foo & echo $!
9322     done
```

9323 The treatment of the signals SIGINT and SIGQUIT with asynchronous lists is described in the  
 9324 Shell and Utilities volume of IEEE Std. 1003.1-200x, Section 2.12, Signals and Error Handling.

9325 Since the connection of the input to the equivalent of /dev/null is considered to occur before  
 9326 redirections, the following script would produce no output:

```
9327     exec < /etc/passwd
9328     cat <&0 &
9329     wait
```

9330 **Sequential Lists**

9331 There is no additional rationale for this section.

9332        **AND Lists**

9333        There is no additional rationale for this section.

9334        **OR Lists**

9335        There is no additional rationale for this section.

9336    *C.2.9.4 Compound Commands*9337        **Grouping Commands**9338        The semicolon shown *{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 2.11, Shell Grammar; <newline> is frequently used.

9341        A proposal was made to use the &lt;do-done&gt; 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.

9346        **For Loop**

9347        The format is shown with generous usage of &lt;newline&gt;s. See the grammar in the Shell and Utilities volume of IEEE Std. 1003.1-200x, Section 2.11, Shell Grammar for a precise description of where &lt;newline&gt;s and semicolons can be interchanged.

9350        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.)9355        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.

9358        PASC Interpretation 1003.2 #169 has been applied changing the grammar.

9359        **Case Conditional Construct**9360        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.9366        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*.9370        The pattern ' \* ', given as the last pattern in a **case** construct, is equivalent to the default case in a C-language **switch** statement.

9371

9372 The grammar shows that reserved words can be used as patterns, even if one is the first word on  
9373 a line. Obviously, the reserved word **esac** cannot be used in this manner.

#### 9374 **If Conditional Construct**

9375 The precise format for the command syntax is described in the Shell and Utilities volume of  
9376 IEEE Std. 1003.1-200x, Section 2.11, Shell Grammar.

#### 9377 **While Loop**

9378 The precise format for the command syntax is described in the Shell and Utilities volume of  
9379 IEEE Std. 1003.1-200x, Section 2.11, Shell Grammar.

#### 9380 **Until Loop**

9381 The precise format for the command syntax is described in the Shell and Utilities volume of  
9382 IEEE Std. 1003.1-200x, Section 2.11, Shell Grammar.

#### 9383 *C.2.9.5 Function Definition Command*

9384 The description of functions in an early proposal was based on the notion that functions should  
9385 behave like miniature shell scripts; that is, except for sharing variables, most elements of an  
9386 execution environment should behave as if they were a new execution environment, and  
9387 changes to these should be local to the function. For example, traps and options should be reset  
9388 on entry to the function, and any changes to them do not affect the traps or options of the caller.  
9389 There were numerous objections to this basic idea, and the opponents asserted that functions  
9390 were intended to be a convenient mechanism for grouping common commands that were to be  
9391 executed in the current execution environment, similar to the execution of the *dot* special built-  
9392 in.

9393 It was also pointed out that the functions described in that early proposal did not provide a local  
9394 scope for everything a new shell script would, such as the current working directory, or *umask*,  
9395 but instead provided a local scope for only a few select properties. The basic argument was that  
9396 if a local scope is needed for the execution environment, the mechanism already existed: the  
9397 application can put the commands in a new shell script and call that script. All historical shells  
9398 that implemented functions, other than the KornShell, have implemented functions that operate  
9399 in the current execution environment. Because of this, traps and options have a global scope  
9400 within a shell script. Local variables within a function were considered and included in another  
9401 early proposal (controlled by the special built-in *local*), but were removed because they do not fit  
9402 the simple model developed for functions and because there was some opposition to adding yet  
9403 another new special built-in that was not part of historical practice. Implementations should  
9404 reserve the identifier *local* (as well as *typeset*, as used in the KornShell) in case this local variable  
9405 mechanism is adopted in a future version of IEEE Std. 1003.1-200x.

9406 A separate issue from the execution environment of a function is the availability of that function  
9407 to child shells. A few objectors maintained that just as a variable can be shared with child shells  
9408 by exporting it, so should a function. In early proposals, the *export* command therefore had a *-f*  
9409 flag for exporting functions. Functions that were exported were to be put into the environment  
9410 as *name()=value* pairs, and upon invocation, the shell would scan the environment for these and  
9411 automatically define these functions. This facility was strongly opposed and was omitted. Some  
9412 of the arguments against exportable functions were as follows:

- 9413 • There was little historical practice. The Ninth Edition shell provided them, but there was  
9414 controversy over how well it worked.



- 9415 • There are numerous security problems associated with functions appearing in the  
9416 environment of a user and overriding standard utilities or the utilities owned by the  
9417 application.
- 9418 • There was controversy over requiring *make* to import functions, where it has historically used  
9419 an *exec* function for many of its command line executions.
- 9420 • Functions can be big and the environment is of a limited size. (The counter-argument was  
9421 that functions are no different than variables in terms of size: there can be big ones, and there  
9422 can be small ones—and just as one does not export huge variables, one does not export huge  
9423 functions. However, this might not apply to the average shell-function writer, who typically  
9424 writes much larger functions than variables.)

9425 As far as can be determined, the functions in the Shell and Utilities volume of  
9426 IEEE Std. 1003.1-200x match those in System V. Earlier versions of the KornShell had two  
9427 methods of defining functions:

```
9428     function fname { compound-list }
```

9429 and:

```
9430     fname() { compound-list }
```

9431 The latter used the same definition as the Shell and Utilities volume of IEEE Std. 1003.1-200x, but  
9432 differed in semantics, as described previously. The current edition of the KornShell aligns the  
9433 latter syntax with the Shell and Utilities volume of IEEE Std. 1003.1-200x and keeps the former as  
9434 is.

9435 The name space for functions is limited to that of a *name* because of historical practice.  
9436 Complications in defining the syntactic rules for the function definition command and in dealing  
9437 with known extensions such as the "@()" usage in the KornShell prevented the name space  
9438 from being widened to a *word*. Using functions to support synonyms such as the "!!" and '%'  
9439 usage in the C shell is thus disallowed to portable applications, but acceptable as an extension.  
9440 For interactive users, the aliasing facilities in the Shell and Utilities volume of  
9441 IEEE Std. 1003.1-200x should be adequate for this purpose. It is recognized that the name space  
9442 for utilities in the file system is wider than that currently supported for functions, if the portable  
9443 file name character set guidelines are ignored, but it did not seem useful to mandate extensions  
9444 in systems for so little benefit to portable applications.

9445 The "()" in the function definition command consists of two operators. Therefore, intermixing  
9446 <blank>s with the *fname*, '( ', and ' )' is allowed, but unnecessary.

9447 An example of how a function definition can be used wherever a simple command is allowed:

```
9448     # If variable i is equal to "yes",
9449     # define function foo to be ls -l
9450     #
9451     [ "$i" = yes ] && foo() {
9452         ls -l
9453     }
```

9454 **C.2.10 Executable Script**

9455 The working group did not reach consensus to adopt this as a core requirement—that is, for  
 9456 POSIX-conforming applications—however, existing practice on UNIX systems indicated that it  
 9457 should be added as an XSI extension, and this was brought into the scope of this revision by The  
 9458 Open Group Base Resolution bwg2000-004. The scope of this feature is to document existing  
 9459 practice and not to invent.

9460 Applications must not assume that the standard utilities will be available in any particular  
 9461 named directory. For example, it cannot be assumed that standard versions of *awk* and *sh* will be  
 9462 available as */bin/sh* or */bin/awk*, respectively, since implementations are permitted to provide  
 9463 non-standard versions of the utilities in these directories.

9464 It is recommended that an installation script for executable scripts use the standard *PATH*  
 9465 returned by a call to the *getconf* utility with the argument *PATH*, combined with the *command*  
 9466 utility to determine the location of a standard utility.

9467 For example, to determine the location of the standard *sh* utility:

```
9468 command -v sh
```

9469 On some systems this might return:

```
9470 /usr/xpg4/bin/sh
```

9471 Note that the installation script should ensure that the returned path name is an absolute path  
 9472 name prior to use, since a shell built-in might be returned for some utilities.

9473 **C.2.11 Shell Grammar**

9474 There are several subtle aspects of this grammar where conventional usage implies rules about  
 9475 the grammar that in fact are not true.

9476 For *compound\_list*, only the forms that end in a *separator* allow a reserved word to be recognized,  
 9477 so usually only a *separator* can be used where a compound list precedes a reserved word (such as  
 9478 **Then**, **Else**, **Do** and **Rbrace**). Explicitly requiring a separator would disallow such valid (if rare)  
 9479 statements as:

```
9480     if (false) then (echo x) else (echo y) fi
```

9481 See the Note under special grammar rule 1.

9482 Concerning the third sentence of rule (1) (“Also, if the parser ...”):

- 9483 • This sentence applies rather narrowly: when a compound list is terminated by some clear  
 9484 delimiter (such as the closing **fi** of an inner **if\_clause**) then it would apply; where the  
 9485 compound list might continue (as in after a *;*), rule (7a) (and consequently the first  
 9486 sentence of rule (1)) would apply. In many instances the two conditions are identical, but this  
 9487 part of rule (1) does not give license to treating a **WORD** as a reserved word unless it is in a  
 9488 place where a reserved word has to appear.

- 9489 • The statement is equivalent to requiring that when the LR(1) lookahead set contains exactly  
 9490 one reserved word, it must be recognized if it is present. (Here “LR(1)” refers to the  
 9491 theoretical concepts, not to any real parser generator.)

9492 For example, in the construct below, and when the parser is at the point marked with ‘^’,  
 9493 the only next legal token is **then** (this follows directly from the grammar rules):

```
9494     if if...fi then ... fi
9495         ^
```

9496 At that point, the **then** must be recognized as a reserved word.

9497 (Depending on the parser generator actually used, “extra” reserved words may be in some  
9498 lookahead sets. It does not really matter if they are recognized, or even if any possible  
9499 reserved word is recognized in that state, because if it is recognized and is not in the  
9500 (theoretical) LR(1) lookahead set, an error is ultimately detected. In the example above, if  
9501 some other reserved word (for example, **while**) is also recognized, an error occurs later.

9502 This is approximately equivalent to saying that reserved words are recognized after other  
9503 reserved words (because it is after a reserved word that this condition occurs), but avoids the  
9504 “except for ...” list that would be required for **case**, **for**, and so on. (Reserved words are of  
9505 course recognized anywhere a *simple\_command* can appear, as well. Other rules take care of  
9506 the special cases of non-recognition, such as rule (4) for **case** statements.)

9507 Note that the body of here-documents are handled by token recognition (see the Shell and  
9508 Utilities volume of IEEE Std. 1003.1-200x, Section 2.3, Token Recognition) and do not appear in  
9509 the grammar directly. (However, the here-document I/O redirection operator is handled as part  
9510 of the grammar.)

9511 The start symbol of the grammar (**complete\_command**) represents either input from the  
9512 command line or a shell script. It is repeatedly applied by the interpreter to its input and  
9513 represents a single “chunk” of that input as seen by the interpreter.

9514 *C.2.11.1 Shell Grammar Lexical Conventions*

9515 There is no additional rationale for this section.

9516 *C.2.11.2 Shell Grammar Rules*

9517 There is no additional rationale for this section.

9518 **C.2.12 Signals and Error Handling**

9519 There is no additional rationale for this section.

9520 **C.2.13 Shell Execution Environment**

9521 Some systems have implemented the last stage of a pipeline in the current environment so that  
9522 commands such as:

9523 `command | read foo`

9524 set variable **foo** in the current environment. This extension is allowed, but not required;  
9525 therefore, a shell programmer should consider a pipeline to be in a subshell environment, but  
9526 not depend on it.

9527 In early proposals, the description of execution environment failed to mention that each  
9528 command in a multiple command pipeline could be in a subshell execution environment. For  
9529 compatibility with some historical shells, the wording was phrased to allow an implementation  
9530 to place any or all commands of a pipeline in the current environment. However, this means that  
9531 a POSIX application must assume each command is in a subshell environment, but not depend  
9532 on it.

9533 The wording about shell scripts is meant to convey the fact that describing “trap actions” can  
9534 only be understood in the context of the shell command language. Outside of this context, such  
9535 as in a C-language program, signals are the operative condition, not traps.

9536 **C.2.14 Pattern Matching Notation**

9537 Pattern matching is a simpler concept and has a simpler syntax than REs, as the former is  
 9538 generally used for the manipulation of file names, which are relatively simple collections of  
 9539 characters, while the latter is generally used to manipulate arbitrary text strings of potentially  
 9540 greater complexity. However, some of the basic concepts are the same, so this section points  
 9541 liberally to the detailed descriptions in the Base Definitions volume of IEEE Std. 1003.1-200x,  
 9542 Chapter 9, Regular Expressions.

9543 *C.2.14.1 Patterns Matching a Single Character*

9544 Both quoting and escaping are described here because pattern matching must work in three  
 9545 separate circumstances:

9546 1. Calling directly upon the shell, such as in path name expansion or in a **case** statement. All  
 9547 of the following match the string or file **abc**:

9548 `abc "abc" a"b" c a\b c a[b] c a["b"] c a[\b] c a["\b"] c a?c a*c`

9549 The following do not:

9550 `"a?c" a*c a\b c`

9551 2. Calling a utility or function without going through a shell, as described for *find* and the  
 9552 *fmitch()* function defined in the System Interfaces volume of IEEE Std. 1003.1-200x.

9553 3. Calling utilities such as *find*, *cpio*, *tar*, or *pax* through the shell command line. In this case,  
 9554 shell quote removal is performed before the utility sees the argument. For example, in:

9555 `find /bin -name "e\c[\h]o" -print`

9556 after quote removal, the backslashes are presented to *find* and it treats them as escape  
 9557 characters. Both precede ordinary characters, so the *c* and *h* represent themselves and *echo*  
 9558 would be found on many historical systems (that have it in **/bin**). To find a file name that  
 9559 contained shell special characters or pattern characters, both quoting and escaping are  
 9560 required, such as:

9561 `pax -r ... "*a\(\?"`

9562 to extract a file name ending with "a(?".

9563 Conforming applications are required to quote or escape the shell special characters (sometimes  
 9564 called metacharacters). If used without this protection, syntax errors can result or  
 9565 implementation extensions can be triggered. For example, the KornShell supports a series of  
 9566 extensions based on parentheses in patterns.

9567 The restriction on a circumflex in a bracket expression is to allow implementations that support  
 9568 pattern matching using the circumflex as the negation character in addition to the exclamation  
 9569 mark. A portable application must use something like "[\^!]" to match either character.

9570 *C.2.14.2 Patterns Matching Multiple Characters*

9571 Since each asterisk matches zero or more occurrences, the patterns "a\*b" and "a\*\*b" have  
 9572 identical functionality.

9573 **Examples**

- 9574 `a[bc]` Matches the strings "ab" and "ac".
- 9575 `a*d` Matches the strings "ad", "abd", and "abcd", but not the string "abc".
- 9576 `a*d*` Matches the strings "ad", "abcd", "abcdef", "aaaad", and "adddd".
- 9577 `*a*d` Matches the strings "ad", "abcd", "efabcd", "aaaad", and "adddd".

9578 **C.2.14.3 Patterns Used for File Name Expansion**

9579 The caveat about a slash within a bracket expression is derived from historical practice. The  
 9580 pattern "a[b/c]d" does not match such path names as **abd** or **a/d**. On some systems (including  
 9581 those conforming to the Single UNIX Specification), it matched a path name of literally  
 9582 "a[b/c]d". On other systems, it produced an undefined condition (an unescaped '[' used  
 9583 outside a bracket expression). In this version, the XSI behavior is now required.

9584 File names beginning with a period historically have been specially protected from view on  
 9585 UNIX systems. A proposal to allow an explicit period in a bracket expression to match a leading  
 9586 period was considered; it is allowed as an implementation extension, but a conforming  
 9587 application cannot make use of it. If this extension becomes popular in the future, it will be  
 9588 considered for a future version of the Shell and Utilities volume of IEEE Std. 1003.1-200x.

9589 Historical systems have varied in their permissions requirements. To match **f\*/bar** has required  
 9590 read permissions on the **f\*** directories in the System V shell, but the Shell and Utilities volume of  
 9591 IEEE Std. 1003.1-200x, the C shell, and KornShell require only search permissions.

9592 **C.2.15 Special Built-In Utilities**

9593 See the RATIONALE sections on the individual reference pages.

**9594 C.3 Batch Environment Services and Utilities****9595 Scope of the Batch Environment Option**

9596 This section summarizes the deliberations of the IEEE P1003.15 (Batch Environment) working  
9597 group in the development of the Batch Environment option, which covers a set of services and  
9598 utilities defining a batch processing system.

9599 This informative section contains historical information concerning the contents of the  
9600 amendment and describes why features were included or discarded by the working group.

**9601 History of Batch Systems**

9602 The supercomputing technical committee began as a “Birds Of a Feather” (BOF) at the January  
9603 1987 Usenix meeting. There was enough general interest to form a supercomputing attachment  
9604 to the /usr/group working groups. Several subgroups rapidly formed. Of those subgroups, the  
9605 batch group was the most ambitious. The first early meetings were spent evaluating user needs  
9606 and existing batch implementations.

9607 To evaluate user needs, individuals from the supercomputing community came and presented  
9608 their needs. Common requests were flexibility, interoperability, control of resources, and ease-  
9609 of-use. Backwards-compatibility was not an issue. The working group then evaluated some  
9610 existing systems. The following different systems were evaluated:

- 9611 • PROD
- 9612 • Convex Distributed Batch
- 9613 • NQS
- 9614 • CTSS
- 9615 • MDQS from Ballistics Research Laboratory (BRL)

9616 Finally, NQS was chosen as a model because it satisfied not only the most user requirements, but  
9617 because it was public domain, already implemented on a variety of hardware platforms, and  
9618 networked-based.

**9619 Historical Implementations of Batch Systems**

9620 Deferred processing of work under the control of a scheduler has been a feature of most  
9621 proprietary operating systems from the earliest days of multi-user systems in order to maximize  
9622 utilization of the computer.

9623 The arrival of UNIX systems proved to be a dilemma to many hardware providers and users  
9624 because it did not include the sophisticated batch facilities offered by the proprietary systems.  
9625 This omission was rectified in 1986 by NASA Ames Research Center who developed the  
9626 Network Queuing System (NQS) as a portable UNIX application that allowed the routing and  
9627 processing of batch “jobs” in a network. To encourage its usage, the product was later put into  
9628 the public domain. It was promptly picked up by UNIX hardware providers, and ported and  
9629 developed for their respective hardware and UNIX implementations.

9630 Many major vendors, who traditionally offer a batch-dominated environment, ported the  
9631 public-domain product to their systems, customized it to support the capabilities of their  
9632 systems, and added many customer-requested features.

9633 Due to the strong hardware provider and customer acceptance of NQS, it was decided to use  
9634 NQS as the basis for the POSIX Batch Environment amendment in 1987. Other batch systems  
9635 considered at the time included CTSS, MDQS (a forerunner of NQS from the Ballistics Research

9636 Laboratory), and PROD (a Los Alamos Labs development). None were thought to have both the  
9637 functionality and acceptability of NQS.

### 9638 **NQS Differences from the at utility**

9639 The base standard *at* and *batch* utilities are not sufficient to meet the batch processing needs in a  
9640 supercomputing environment and additional functionality in the areas of resource management,  
9641 job scheduling, system management, and control of output is required.

### 9642 **Batch Environment Option Definitions**

9643 The concept of a batch job is closely related to a session with a session leader. The main  
9644 difference is that a batch job does not have a controlling terminal. There has been much debate  
9645 over whether to use the term *request* or *job*. Job was the final choice because of the historical use  
9646 of this term in the batch environment.

9647 The current definition for job identifiers is not sufficient with the model of destinations. The  
9648 current definition is:

9649 `sequence_number.originating_host`

9650 Using the model of destination, a host may include multiple batch nodes, the location of which is  
9651 identified uniquely by a name or directory service. If the current definition is used, batch nodes  
9652 running on the same host would have to coordinate their use of sequence numbers, as sequence  
9653 numbers are assigned by the originating host. The alternative is to use the originating batch node  
9654 name instead of the originating host name.

9655 The reasons for wishing to run more than one batch system per host could be the following:

9656 A test and production batch system are maintained on a single host. This is most likely in a  
9657 development facility, but could also arise when a site is moving from one version to another.  
9658 The new batch system could be installed as a test version that is completely separate from the  
9659 production batch system, so that problems can be isolated to the test system. Requiring the batch  
9660 nodes to coordinate their use of sequence numbers creates a dependency between the two  
9661 nodes, and that defeats the purpose of running two nodes.

9662 A site has multiple departments using a single host, with different management policies. An  
9663 example of contention might be in job selection algorithms. One group might want a FIFO type  
9664 of selection, while another group wishes to use a more complex algorithm based on resource  
9665 availability. Again, requiring the batch nodes to coordinate is an unnecessary binding.

9666 The proposal eventually accepted was to replace originating host with originating batch node.  
9667 This supplies sufficient granularity to ensure unique job identifiers. If more than one batch node  
9668 is on a particular host, they each have their own unique name.

9669 The queue portion of a destination is not part of the job identifier as these are not required to be  
9670 unique between batch nodes. For instance, two batch nodes may both have queues called small,  
9671 medium, and large. It is only the batch node name that is uniquely identifiable throughout the  
9672 batch system. The queue name has no additional function in this context.

9673 Assume there are three batch nodes, each of which has its own name server. On batch node one,  
9674 there are no queues. On batch node two, there are fifty queues. On batch node three, there are  
9675 forty queues. The system administrator for batch node one does not have to configure queues,  
9676 because there are none implemented. However, if a user wishes to send a job to either batch  
9677 node two or three, the system administrator for batch node one must configure a destination  
9678 that maps to the appropriate batch node and queue. If every queue is to be made accessible from  
9679 batch node one, the system administrator has to configure ninety destinations.

9680 To avoid requiring this, there should be a mechanism to allow a user to separate the destination  
9681 into a batch node name and a queue name. Then, an implementation that is configured to get to  
9682 all the batch nodes does not need any more configuration to allow a user to get to all of the  
9683 queues on all of the batch nodes. The node name is used to locate the batch node, while the  
9684 queue name is sent unchanged to that batch node.

9685 The following are requirements that a destination identifier must be capable of providing:

- 9686 • The ability to direct a job to a queue in a particular batch node.
- 9687 • The ability to direct a job to a particular batch node.
- 9688 • The ability to group at a higher level than just one queue. This includes grouping similar  
9689 queues across multiple batch nodes (this is a pipe queue today).
- 9690 • The ability to group batch nodes. This allows a user to submit a job to a group name with no  
9691 knowledge of the batch node configuration. This also provides aliasing as a special case.  
9692 Aliasing is a group containing only one batch node name. The group name is the alias.

9693 In addition, the administrator has the following requirements:

- 9694 • The ability to control access to the queues.
- 9695 • The ability to control access to the batch nodes.
- 9696 • The ability to control access to groups of queues (pipe queues).
- 9697 • The ability to configure retry time intervals and durations.

9698 The requirements of the user are met by destination as explained in the following:

9699 The user has the ability to specify a queue name, which is known only to the batch node  
9700 specified. There is no configuration of these queues required on the submitting node.

9701 The user has the ability to specify a batch node whose name is network-unique. The  
9702 configuration required is that the batch node be defined as an application, just as other  
9703 applications such as FTP are configured.

9704 Once a job reaches a queue, it can again become a user of the batch system. The batch node can  
9705 choose to send the job to another batch node or queue or both. In other words, the routing is at  
9706 an application level, and it is up to the batch system to choose where the job will be sent.  
9707 Configuration is up to the batch node where the queue resides. This provides grouping of  
9708 queues across batch nodes or within a batch node. The user submits the job to a queue, which by  
9709 definition routes the job to other queues or nodes or both.

9710 A node name may be given to a naming service, which returns multiple addresses as opposed to  
9711 just one. This provides grouping at a batch node level. This is a local issue, meaning that the  
9712 batch node must choose only one of these addresses. The list of addresses is not sent with the  
9713 job, and once the job is accepted on another node, there is no connection between the list and the  
9714 job. The requirements of the administrator are met by destination as explained in the following:

9715 The control of queues is a batch system issue, and will be done using the batch administrative  
9716 utilities.

9717 The control of nodes is a network issue, and will be done through whatever network facilities  
9718 are available.

9719 The control of access to groups of queues (pipe queues) is covered by the control of any other  
9720 queue. The fact that the job may then be sent to another destination is not relevant.

9721 The propagation of a job across more than one point-to-point connection was dropped because  
9722 of its complexity and because all of the issues arising from this capability could not be resolved.



9723 It could be provided as additional functionality at some time in the future.

9724 The addition of *network* as a defined term was done to clarify the difference between a network  
9725 of batch nodes as opposed to a network of hosts. A network of batch nodes is referred to as a  
9726 batch system. The network refers to the actual host configuration. A single host may have  
9727 multiple batch nodes.

9728 In the absence of a standard network naming convention, this option establishes its own  
9729 convention for the sake of consistency and expediency. This is subject to change, should a future  
9730 working group develop a standard naming convention for network path names.

### 9731 **C.3.1 Batch General Concepts**

9732 During the development of the Batch Environment option, a number of topics were discussed at  
9733 length which influenced the wording of the normative text but could not be included in the final  
9734 text. The following items are some of the most significant terms and concepts of those discussed:

- 9735 • Small and Consistent Command Set

9736 Often, conventional utilities from UNIX systems have a very complicated utility syntax and  
9737 usage. This can often result in confusion and errors when trying to use them. The Batch  
9738 Environment option utility set, on the other hand, has been paired to a small set of robust  
9739 utilities with an orthogonal calling sequence.

- 9740 • Checkpoint/Restart

9741 This feature permits an already executing process to checkpoint or save its contents. Some  
9742 implementations permit this at both the batch utility level; for example, checkpointing this  
9743 job upon its abnormal termination or from within the job itself via a system call. Support of  
9744 checkpoint/restart is optional. A conscious, careful effort was made to make the *qsub* and  
9745 *qmgr* utilities consistently refer to checkpoint/restart as optional functionality.

- 9746 • Rerunability

9747 When a user submits a job for batch processing, they can designate it “rerunnable” in that it  
9748 will automatically resume execution from the start of the job if the machine on which it was  
9749 executing crashes for some reason. The decision on whether the job will be rerun or not is  
9750 entirely up to the submitter of the job and no decisions will be made within the batch system.  
9751 A job that is rerunnable and has been submitted with the proper checkpoint/restart switch  
9752 will first be checkpointed and execution begun from that point. Furthermore, use of the  
9753 implementation-defined checkpoint/restart feature will be not be defined in this context.

- 9754 • Error Codes

9755 All utilities exit with error status zero (0) if successful, one (1) if a user error occurred, and  
9756 two (2) for an internal Batch Environment option error.

- 9757 • Level of Portability

9758 Portability is specified at both the user, operator, and administrator levels. A conforming  
9759 batch implementation prevents identical functionality and behavior at all these levels.  
9760 Additionally, portable batch shell scripts with embedded Batch Environment option utilities  
9761 adds an additional level of portability.

- 9762 • Resource Specification

9763 A small set of globally understood resources, such as memory and CPU time, is specified. All  
9764 conforming batch implementations are able to process them in a manner consistent with the  
9765 yet-to-be-developed resource management model. Resources not in this amendment set are  
9766 ignored and passed along as part of the argument stream of the utility.

- 9767
  - Maximum of 80 Characters on Output Display
- 9768

At one time, existing displays were limited to 80 characters in length for purposes of
- 9769

readability, both in the amendment and online. Current internationalization efforts
- 9770

discourage specifying displays in the normative text. Instead, all suggested displays appear
- 9771

in an informative annex for illustrative purposes. As before, length is limited to 80 characters
- 9772

for readability purposes.
- 9773
  - Queue Position
- 9774

Queue position is the place a job occupies in a queue. It is dependent on a variety of factors
- 9775

such as submission time and priority. Since priority may be affected by the implementation
- 9776

of fair share scheduling, the definition of queue position is implementation-defined.
- 9777
  - Queue ID
- 9778

A numerical queue ID is an external requirement for purposes of accounting. The
- 9779

identification number was chosen over queue name for processing convenience.
- 9780
  - Job ID
- 9781

A common notion of “jobs” is a collection of processes whose process group cannot be
- 9782

altered and is used for resource management and accounting. This concept is
- 9783

implementation-defined and, as such, has been omitted from the batch amendment.
- 9784
  - Bytes versus Words
- 9785

Except for one case, bytes are used as the standard unit for memory size. Furthermore, the
- 9786

definition of a word varies from machine to machine. Therefore, bytes will be the default unit
- 9787

of memory size.
- 9788
  - Regular Expressions
- 9789

The standard definition of regular expressions is much too broad to be used in the batch
- 9790

utility syntax. All that is needed is a simple concept of “all”; for example, delete all my jobs
- 9791

from the named queue. For this reason, regular expressions have been eliminated from the
- 9792

batch amendment.
- 9793
  - Display Privacy
- 9794

How much data should be displayed locally through library functions? Local policy dictates
- 9795

the amount of privacy. Library functions must be used to create and enforce local policy.
- 9796

Network and local *qstats* must reflect the policy of the server machine.
- 9797
  - Remote Host Naming Convention
- 9798

It was decided that host names would be a maximum of 255 characters in length, with at
- 9799

most 15 characters being shown in displays. The 255 character limit was chosen because it is
- 9800

consistent with BSD). The 15-character limit was an arbitrary decision.
- 9801
  - Network Administration
- 9802

Network administration is important, but is outside the scope of the batch amendment.
- 9803

Network administration could done with *rsh*. However, authentication becomes two-sided.
- 9804
  - Network Administration Philosophy
- 9805

Keep it simple. Centralized management should be possible. For example, Los Alamos needs
- 9806

a dumb set of CPUs to be managed by a central system *versus* several independently-
- 9807

managed systems as is the general case for the Batch Environment option.

- 9808       • Operator Utility Defaults (that is, Default Host, User, Account, and so on)
- 9809        It was decided that usability would override orthogonality and syntactic consistency.
- 9810       • The Batch System Manager and Operator Distinction
- 9811        The distinction between manager and operator is that operators can only control the flow of
- 9812        jobs. A manager can alter the batch system configuration in addition to job flow. POSIX
- 9813        makes a distinction between user and system administrator but goes no further. The
- 9814        concepts of manager and operator privileges fall under local policy. The distinction between
- 9815        manager and operator is historical in batch environments, and the Batch Environment option
- 9816        has continued that distinction.
- 9817       • The Batch System Administrator
- 9818        An administrator is equivalent to a batch system manager.
- 9819       • Network Administration *versus* Checkpoint/Recovery
- 9820        Network administration is better left for a future revision of IEEE Std. 1003.1-200x. If
- 9821        network administration is put up against checkpoint/recovery, the argument is that there are
- 9822        possible solutions for network administration. However, checkpoint/recovery currently has
- 9823        no solution. This is another reason the issue of checkpoint/recovery should be addressed
- 9824        first.

### 9825 **C.3.2 Batch Services**

- 9826       This rationale is provided as informative rather than normative text, to avoid placing
- 9827       requirements on implementors regarding the use of symbolic constants, but at the same time to
- 9828       give implementors a preferred practice for assigning values to these constants to promote
- 9829       interoperability.
- 9830       The *Checkpoint* and *Minimum\_Cpu\_Interval* attributes induce a variety of behavior depending
- 9831       upon their values. Some jobs cannot or should not be checkpointed. Other users will simply
- 9832       need to ensure job continuation across planned downtimes; for example, scheduled preventive
- 9833       maintenance. For users consuming expensive resources, or for jobs that run longer than the
- 9834       mean time between failures, however, periodic checkpointing may be essential. However,
- 9835       system administrators must be able to set minimum checkpoint intervals on a queue-by-queue
- 9836       basis to guard against; for example, naive users specifying interval values too small on memory
- 9837       intensive jobs. Otherwise, system overhead would adversely affect performance.
- 9838       The use of symbolic constants, such as `NO_CHECKPOINT`, was introduced to lend a degree of
- 9839       formalism and portability to this option.
- 9840       Support for checkpointing is optional for servers. However, clients must provide for the `-c`
- 9841       option, since in a distributed environment the job may run on a server that does provide such
- 9842       support, even if the host of the client does not support the checkpoint feature.
- 9843       If the user does not specify the `-c` option, the default action is left unspecified by this option.
- 9844       Some implementations may wish to do checkpointing by default; others may wish to checkpoint
- 9845       only under an explicit request from the user.
- 9846       The *Priority* attribute has been made non-optional. All clients already had been required to
- 9847       support the `-p` option. The concept of prioritization is common in historical implementations.
- 9848       The default priority is left to the server to establish.
- 9849       The *Hold\_Types* attribute has been modified to allow for implementation-defined hold types to
- 9850       be passed to a batch server.

9851 It was the intent of the IEEE P1003.15 working group to mandate the support for the  
9852 *Resource\_List* attribute in this option by referring to another amendment, specifically P1003.1a.  
9853 However, during the development of P1003.1a this was excluded. As such this requirement has  
9854 been removed from the normative text.

9855 The *Shell\_Path* attribute has been modified to accept a list of shell paths that are associated with  
9856 a host. The name of the attribute has been changed to *Shell\_Path\_List*.

9857 **C.3.3 Common Behavior for Batch Environment Utilities**

9858 This section was defined to meet the goal of a “Small and Consistent Command Set” for this  
9859 option.

9860 **C.4 Utilities**

9861 See the RATIONALE sections on the individual reference pages.



9862 / *Rationale (Informative)*

9863 **Part D:**

9864 **Portability Considerations**

9865 *The Open Group*





## Portability Considerations (Informative)

9866

9867 This section contains information to satisfy various international requirements:

- 9868 • Section D.1 describes perceived user requirements.
- 9869 • Section D.2 (on page 3560) indicates how the facilities of IEEE Std. 1003.1-200x satisfy those  
9870 requirements.
- 9871 • Section D.3 (on page 3567) offers guidance to writers of profiles on how the configurable  
9872 options, limits, and optional behavior of IEEE Std. 1003.1-200x should be cited in profiles.

### 9873 D.1 User Requirements

9874 This section describes the user requirements that were perceived by the developers of  
9875 IEEE Std. 1003.1-200x. The primary source for these requirements was an analysis of historical  
9876 practice in widespread use, as typified by the base documents listed in Section A.1.1 (on page  
9877 3311).

9878 IEEE Std. 1003.1-200x addresses the needs of users requiring open systems solutions for source  
9879 code portability of applications. It currently addresses users requiring open systems solutions  
9880 for source-code portability of applications involving multi-programming and process  
9881 management (creating processes, signaling, and so on); access to files and directories in a  
9882 hierarchy of file systems (opening, reading, writing, deleting files, and so on); access to  
9883 asynchronous communications ports and other special devices; access to information about  
9884 other users of the system; facilities supporting applications requiring bounded (realtime)  
9885 response.

9886 The following users are identified for IEEE Std. 1003.1-200x:

- 9887 • Those employing applications written in high-level languages, such as C, Ada, or FORTRAN.
- 9888 • Users who desire portable applications that do not necessarily require the characteristics of  
9889 high-level languages (for example, the speed of execution of compiled languages or the  
9890 relative security of source code intellectual property inherent in the compilation process).
- 9891 • Users who desire portable applications that can be developed quickly and can be modified  
9892 readily without the use of compilers and other system components that may be unavailable  
9893 on small systems or those without special application development capabilities.
- 9894 • Users who interact with a system to achieve general-purpose time-sharing capabilities  
9895 common to most business or government offices or academic environments: editing, filing,  
9896 inter-user communications, printing, and so on.
- 9897 • Users who develop applications for POSIX-conformant systems.
- 9898 • Users who develop applications for UNIX systems.

9899 An acknowledged restriction on applicable users is that they are limited to the group of  
9900 individuals who are familiar with the style of interaction characteristic of historically-derived  
9901 systems based on one of the UNIX operating systems (as opposed to other historical systems  
9902 with different models, such as MS/DOS, Macintosh, VMS, MVS, and so on). Typical users  
9903 would include program developers, engineers, or general-purpose time-sharing users.

9904 The requirements of users of IEEE Std. 1003.1-200x can be summarized as a single goal:  
9905 *application source portability*. The requirements of the user are stated in terms of the requirements  
9906 of portability of applications. This in turn becomes a requirement for a standardized set of  
9907 syntax and semantics for operations commonly found on many operating systems.

9908 The following sections list the perceived requirements for application portability.

#### 9909 **D.1.1 Configuration Interrogation**

9910 An application must be able to determine whether and how certain optional features are  
9911 provided and to identify the system upon which it is running, so that it may appropriately adapt  
9912 to its environment.

9913 Applications must have sufficient information to adapt to varying behaviors of the system.

#### 9914 **D.1.2 Process Management**

9915 An application must be able to manage itself, either as a single process or as multiple processes.  
9916 Applications must be able to manage other processes when appropriate.

9917 Applications must be able to identify, control, create, and delete processes, and there must be  
9918 communication of information between processes and to and from the system.

9919 Applications must be able to use multiple flows of control with a process (threads) and  
9920 synchronize operations between these flows of control.

#### 9921 **D.1.3 Access to Data**

9922 Applications must be able to operate on the data stored on the system, access it, and transmit it  
9923 to other applications. Information must have protection from unauthorized or accidental access  
9924 or modification.

#### 9925 **D.1.4 Access to the Environment**

9926 Applications must be able to access the external environment to communicate their input and  
9927 results.

#### 9928 **D.1.5 Access to Determinism and Performance Enhancements**

9929 Applications must have sufficient control of resource allocation to ensure the timeliness of  
9930 interactions with external objects.

#### 9931 **D.1.6 Operating System-Dependent Profile**

9932 The capabilities of the operating system may make certain optional characteristics of the base  
9933 language in effect no longer optional, and this should be specified.

**9934 D.1.7 I/O Interaction**

9935 The interaction between the C language I/O subsystem (*stdio*) and the I/O subsystem of  
9936 IEEE Std. 1003.1-200x must be specified.

**9937 D.1.8 Internationalization Interaction**

9938 The effects of the environment of IEEE Std. 1003.1-200x on the internationalization facilities of  
9939 the C language must be specified.

**9940 D.1.9 C-Language Extensions**

9941 Certain functions in the C language must be extended to support the additional capabilities  
9942 provided by IEEE Std. 1003.1-200x.

**9943 D.1.10 Command Language**

9944 Users should be able to define procedures that combine simple tools and/or applications into  
9945 higher-level components that perform to the specific needs of the user. The user should be able  
9946 to store, recall, use, and modify these procedures. These procedures should employ a powerful  
9947 command language that is used for recurring tasks in portable applications (scripts) in the same  
9948 way that it is used interactively to accomplish one-time tasks. The language and the utilities that  
9949 it uses must be consistent between systems to reduce errors and retraining.

**9950 D.1.11 Interactive Facilities**

9951 Use the system to accomplish individual tasks at an interactive terminal. The interface should be  
9952 consistent, intuitive, and offer usability enhancements to increase the productivity of terminal  
9953 users, reduce errors, and minimize retraining costs. Online documentation or usage assistance  
9954 should be available.

**9955 D.1.12 Accomplish Multiple Tasks Simultaneously**

9956 Access applications and interactive facilities from a single terminal without requiring serial  
9957 execution: switch between multiple interactive tasks; schedule one-time or periodic background  
9958 work; display the status of all work in progress or scheduled; influence the priority scheduling of  
9959 work, when authorized.

**9960 D.1.13 Complex Data Manipulation**

9961 Manipulate data in files in complex ways: sort, merge, compare, translate, edit, format, pattern  
9962 match, select subsets (strings, columns, fields, rows, and so on). These facilities should be  
9963 available to both portable applications and interactive users.

**9964 D.1.14 File Hierarchy Manipulation**

9965 Create, delete, move/rename, copy, backup/archive, and display files and directories. These  
9966 facilities should be available to both portable applications and interactive users.

**9967 D.1.15 Locale Configuration**

9968 Customize applications and interactive sessions for the cultural and language conventions of the  
9969 user. Employ a wide variety of standard character encodings. These facilities should be available  
9970 to both portable applications and interactive users.

**9971 D.1.16 Inter-User Communication**

9972 Send messages or transfer files to other users on the same system or other systems on a network.  
9973 These facilities should be available to both portable applications and interactive users.

**9974 D.1.17 System Environment**

9975 Display information about the status of the system (activities of users and their interactive and  
9976 background work, file system utilization, system time, configuration, and presence of optional  
9977 facilities) and the environment of the user (terminal characteristics, and so on). Inform the  
9978 system operator/administrator of problems. Control access to user files and other resources.

**9979 D.1.18 Printing**

9980 Output files on a variety of output device classes, accessing devices on local or network-  
9981 connected systems. Control (or influence) the formatting, priority scheduling, and output  
9982 distribution of work. These facilities should be available to both portable applications and  
9983 interactive users.

**9984 D.1.19 Software Development**

9985 Develop (create and manage source files, compile/interpret, debug) portable open systems  
9986 applications and package them for distribution to, and updating of, other systems.

**9987 D.2 Portability Capabilities**

9988 This section describes the significant portability capabilities of IEEE Std. 1003.1-200x and  
9989 indicates how the user requirements listed in Section D.1 (on page 3557) are addressed. The  
9990 capabilities are listed in the same format as the preceding user requirements; they are  
9991 summarized below:

- 9992 • Configuration Interrogation
- 9993 • Process Management
- 9994 • Access to Data
- 9995 • Access to the Environment
- 9996 • Access to Determinism and Performance Enhancements
- 9997 • Operating System-Dependent Profile
- 9998 • I/O Interaction
- 9999 • Internationalization Interaction
- 10000 • C-Language Extensions
- 10001 • Command Language
- 10002 • Interactive Facilities

- 10003 • Accomplish Multiple Tasks Simultaneously
- 10004 • Complex Data Manipulation
- 10005 • File Hierarchy Manipulation
- 10006 • Locale Configuration
- 10007 • Inter-User Communication
- 10008 • System Environment
- 10009 • Printing
- 10010 • Software Development

### 10011 **D.2.1 Configuration Interrogation**

10012 The *uname()* operation provides basic identification of the system. The *sysconf()*, *pathconf()*, and  
 10013 *fpathconf()* functions and the *getconf* utility provide means to interrogate the implementation to  
 10014 determine how to adapt to the environment in which it is running. These values can be either  
 10015 static (indicating that all instances of the implementation have the same value) or dynamic  
 10016 (indicating that different instances of the implementation have the different values, or that the  
 10017 value may vary for other reasons, such as reconfiguration).

### 10018 **Unsatisfied Requirements**

10019 None directly. However, as new areas are added, there will be a need for additional capability in  
 10020 this area.

### 10021 **D.2.2 Process Management**

10022 The *fork()*, *exec* family, and *spawn()* functions provide for the creation of new processes or the  
 10023 insertion of new applications into existing processes. The *\_Exit()*, *\_exit()*, *exit()*, and *abort()*  
 10024 functions allow for the termination of a process by itself. The *wait()* and *waitpid()* functions  
 10025 allow one process to deal with the termination of another.

10026 The *times()* function allows for basic measurement of times used by a process. Various  
 10027 functions, including *fstat()*, *getegid()*, *geteuid()*, *getgid()*, *getrgid()*, *getgrnam()*, *getlogin()*,  
 10028 *getpid()*, *getppid()*, *getpwnam()*, *getpwuid()*, *getuid()*, *lstat()*, and *stat()*, provide for access to the  
 10029 identifiers of processes and the identifiers and names of owners of processes (and files).

10030 The various functions operating on environment variables provide for communication of  
 10031 information (primarily user-configurable defaults) from a parent to child processes.

10032 The operations on the current working directory control and interrogate the directory from  
 10033 which relative file name searches start. The *umask()* function controls the default protections  
 10034 applied to files created by the process.

10035 The *alarm()*, *pause()*, *sleep()*, *ualarm()*, and *usleep()* operations allow the process to suspend until  
 10036 a timer has expired or to be notified when a period of time has elapsed. The *time()* operation  
 10037 interrogates the current time and date.

10038 The signal mechanism provides for communication of events either from other processes or  
 10039 from the environment to the application, and the means for the application to control the effect  
 10040 of these events. The mechanism provides for external termination of a process and for a process  
 10041 to suspend until an event occurs. The mechanism also provides for a value to be associated with  
 10042 an event.

10043 Job control provides a means to group processes and control them as groups, and to control their  
10044 access to the function between the user and the system (the *controlling terminal*). It also provides  
10045 the means to suspend and resume processes.

10046 The Process Scheduling option provides control of the scheduling and priority of a process.

10047 The Message Passing option provides a means for interprocess communication involving small  
10048 amounts of data.

10049 The Memory Management facilities provide control of memory resources and for the sharing of  
10050 memory. This functionality is mandatory on XSI-conformant systems.

10051 The Threads facilities provide multiple flows of control with a process (threads),  
10052 synchronization between threads, association of data with threads, and controlled cancelation of  
10053 threads.

10054 The XSI interprocess communications functionality provide an alternate set of facilities to  
10055 manipulate semaphores, message queues, and shared memory. These are provided on XSI-  
10056 conformant systems to support portable applications developed to run on UNIX systems.

### 10057 **D.2.3 Access to Data**

10058 The *open()*, *close()*, *fclose()*, *fopen()*, and *pipe()* functions provide for access to files and data.  
10059 Such files may be regular files, interprocess data channels (pipes), or devices. Additional types  
10060 of objects in the file system are permitted and are being contemplated for standardization.

10061 The *access()*, *chmod()*, *chown()*, *dup()*, *dup2()*, *fchmod()*, *fcntl()*, *fstat()*, *ftruncate()*, *lstat()*,  
10062 *readlink()*, *realpath()*, *stat()*, and *utime()* functions allow for control and interrogation of file and  
10063 file-related objects, (including symbolic links) and their ownership, protections, and timestamps.

10064 The *fgetc()*, *fputc()*, *fread()*, *fseek()*, *fsetpos()*, *fwrite()*, *getc()*, *getch()*, *lseek()*, *putchar()*, *putc()*,  
10065 *read()*, and *write()* functions provide for data transfer from the application to files (in all their  
10066 forms).

10067 The *closedir()*, *link()*, *mkdir()*, *opendir()*, *readdir()*, *rename()*, *rmdir()*, *rewinddir()*, and *unlink()*  
10068 functions provide for a complete set of operations on directories. Directories can arbitrarily  
10069 contain other directories, and a single file can be mentioned in more than one directory.

10070 The file-locking mechanism provides for advisory locking (protection during transactions) of  
10071 ranges of bytes (in effect, records) in a file.

10072 The *confstr()*, *fpathconf()*, *pathconf()*, and *sysconf()* functions provide for enquiry as to the  
10073 behavior of the system where variability is permitted.

10074 The Synchronized Input and Output option provides for assured commitment of data to media.

10075 The Asynchronous Input and Output option provides for initiation and control of asynchronous  
10076 data transfers.

### 10077 **D.2.4 Access to the Environment**

10078 The operations and types in the Base Definitions volume of IEEE Std. 1003.1-200x, Chapter 11,  
10079 General Terminal Interface are provided for access to asynchronous serial devices. The primary  
10080 intended use for these is the controlling terminal for the application (the interaction point  
10081 between the user and the system). They are general enough to be used to control any  
10082 asynchronous serial device. The functions are also general enough to be used with many other  
10083 device types as a user interface when some emulation is provided.

10084 Less detailed access is provided for other device types, but in many instances an application  
10085 need not know whether an object in the file system is a device or a regular file to operate  
10086 correctly.

10087 **Unsatisfied Requirements**

10088 Detailed control of common device classes, specifically magnetic tape, is not provided.

10089 **D.2.5 Bounded (Realtime) Response**

10090 The Realtime Signals Extension provides queued signals and the prioritization of the handling of  
10091 signals. The SCHED\_FIFO, SCHED\_SPORADIC, and SCHED\_RR scheduling policies provide  
10092 control over processor allocation. The Semaphores option provides high-performance  
10093 synchronization. The Memory Management functions provide memory locking for control of  
10094 memory allocation, file mapping for high-performance, and shared memory for high-  
10095 performance interprocess communication. The Message Passing option provides for interprocess  
10096 communication without being dependent on shared memory.

10097 The Timers option provides a high resolution function called *nanosleep()* with a finer resolution  
10098 than the *sleep()* function.

10099 The Typed Memory Objects option, the Monotonic Clock option, and the Timeouts option  
10100 provide further facilities for applications to use to obtain predictable bounded response.

10101 **D.2.6 Operating System-Dependent Profile**

10102 IEEE Std. 1003.1-200x makes no distinction between text and binary files. The values of  
10103 EXIT\_SUCCESS and EXIT\_FAILURE are further defined.

10104 **Unsatisfied Requirements**

10105 None known, but the ISO C standard may contain some additional options that could be  
10106 specified.

10107 **D.2.7 I/O Interaction**

10108 IEEE Std. 1003.1-200x defines how each of the ISO C standard *stdio* functions interact with the  
10109 POSIX.1 operations, typically specifying the behavior in terms of POSIX.1 operations.

10110 **Unsatisfied Requirements**

10111 None.

10112 **D.2.8 Internationalization Interaction**

10113 The IEEE Std. 1003.1-200x environment operations provide a means to define the environment  
10114 for *setlocale()* and time functions such as *ctime()*. The *tzset()* function is provided to set time  
10115 conversion information.

10116 The *nl\_langinfo()* function is provided as an XSI extension to query locale-specific cultural  
10117 settings.

10118 **Unsatisfied Requirements**

10119 None.

10120 **D.2.9 C-Language Extensions**10121 The *setjmp()* and *longjmp()* functions are not defined to be cognizant of the signal masks defined  
10122 for POSIX.1. The *sigsetjmp()* and *siglongjmp()* functions are provided to fill this gap.10123 The *\_setjmp()* and *\_longjmp()* functions are provided as XSI extensions to support historic  
10124 practice.10125 **Unsatisfied Requirements**

10126 None.

10127 **D.2.10 Command Language**10128 The shell command language, as described in Shell and Utilities volume of  
10129 IEEE Std. 1003.1-200x, Chapter 2, Shell Command Language, is a common language useful in  
10130 batch scripts, through an API to high-level languages (for the C-Language Binding option,  
10131 *system()* and *popen()*) and through an interactive terminal (see the *sh* utility). The shell language  
10132 has many of the characteristics of a high-level language, but it has been designed to be more  
10133 suitable for user terminal entry and includes interactive debugging facilities. Through the use of  
10134 pipelining, many complex commands can be constructed from combinations of data filters and  
10135 other common components. Shell scripts can be created, stored, recalled, and modified by the  
10136 user with simple editors.10137 In addition to the basic shell language, the following utilities offer features that simplify and  
10138 enhance programmatic access to the utilities and provide features normally found only in high-  
10139 level languages: *basename*, *bc*, *command*, *dirname*, *echo*, *env*, *expr*, *false*, *printf*, *read*, *sleep*, *tee*, *test*,  
10140 *time\**,<sup>2</sup> *true*, *wait*, *xargs*, and all of the special built-in utilities in the Shell and Utilities volume of  
10141 IEEE Std. 1003.1-200x, Section 2.15, Special Built-In Utilities .10142 **Unsatisfied Requirements**

10143 None.

10144 **D.2.11 Interactive Facilities**10145 The utilities offer a common style of command-line interface through conformance to the Utility  
10146 Syntax Guidelines (see the Base Definitions volume of IEEE Std. 1003.1-200x, Section 12.2, Utility  
10147 Syntax Guidelines) and the common utility defaults (see the Shell and Utilities volume of  
10148 IEEE Std. 1003.1-200x, Section 1.11, Utility Description Defaults). The *sh* utility offers an  
10149 interactive command-line history and editing facility. The following utilities in the User  
10150 Portability Utilities option have been customized for interactive use: *alias*, *ex*, *fc*, *mailx*, *more*, *talk*,  
10151 *vi*, *unalias*, and *write*; the *man* utility offers online access to system documentation.

10152 \_\_\_\_\_

10153 2. The utilities listed with an asterisk here and later in this section are present only on systems which support the User Portability  
10154 Utilities option. There may be further restrictions on the utilities offered with various configuration option combinations; see the  
10155 individual utility descriptions.



10156 **Unsatisfied Requirements**

10157 The command line interface to individual utilities is as intuitive and consistent as historical  
 10158 practice allows. Work underway based on graphical user interfaces may be more suitable for  
 10159 novice or occasional users of the system.

10160 **D.2.12 Accomplish Multiple Tasks Simultaneously**

10161 The shell command language offers background processing through the asynchronous list  
 10162 command form; see the Shell and Utilities volume of IEEE Std. 1003.1-200x, Section 2.9, Shell  
 10163 Commands. The *nohup* utility makes background processing more robust and usable. The *kill*  
 10164 utility can terminate background jobs. When the User Portability Utilities option is supported,  
 10165 the following utilities allow manipulation of jobs: *bg*, *fg*, and *jobs*. Also, if the User Portability  
 10166 Utilities option is supported, the following can support periodic job scheduling, control, and  
 10167 display: *at*, *batch*, *crontab*, *nice*, *ps*, and *renice*.

10168 **Unsatisfied Requirements**

10169 Terminals with multiple windows may be more suitable for some multi-tasking interactive uses  
 10170 than the job control approach in IEEE Std. 1003.1-200x. See the comments on graphical user  
 10171 interfaces in Section D.2.11 (on page 3564). The *nice* and *renice* utilities do not necessarily take  
 10172 advantage of complex system scheduling algorithms that are supported by the realtime options  
 10173 within IEEE Std. 1003.1-200x.

10174 **D.2.13 Complex Data Manipulation**

10175 The following utilities address user requirements in this area: *asa*, *awk*, *bc*, *cmp*, *comm*, *csplit\**, *cut*,  
 10176 *dd*, *diff*, *ed*, *ex\**, *expand\**, *expr*, *find*, *fold*, *grep*, *head*, *join*, *od*, *paste*, *pr*, *printf*, *sed*, *sort*, *split\**, *tabs\**, *tail*,  
 10177 *tr*, *unexpand\**, *uniq*, *uudecode\**, *uuencode\**, and *wc*.

10178 **Unsatisfied Requirements**

10179 Sophisticated text formatting utilities, such as *troff* or *TeX*, are not included. Standards work in  
 10180 the area of SGML may satisfy this.

10181 **D.2.14 File Hierarchy Manipulation**

10182 The following utilities address user requirements in this area: *basename*, *cd*, *chgrp*, *chmod*, *chown*,  
 10183 *cksum*, *cp*, *dd*, *df\**, *diff*, *dirname*, *du\**, *find*, *ls*, *ln*, *mkdir*, *mkfifo*, *mv*, *patch\**, *pathchk*, *pax*, *pwd*, *rm*, *rmdir*,  
 10184 *test*, and *touch*.

10185 **Unsatisfied Requirements**

10186 Some graphical user interfaces offer more intuitive file manager components that allow file  
 10187 manipulation through the use of icons for novice users.

**10188 D.2.15 Locale Configuration**

10189 The standard utilities are affected by the various *LC\_* variables to achieve locale-dependent  
10190 operation: character classification, collation sequences, regular expressions and shell pattern  
10191 matching, date and time formats, numeric formatting, and monetary formatting. When the  
10192 POSIX2\_LOCALEDEF option is supported, applications can provide their own locale definition  
10193 files. The following utilities address user requirements in this area: *date*, *ed*, *ex\**, *find*, *grep*, *locale*,  
10194 *localedef*, *more\**, *sed*, *sh*, *sort*, *tr*, *uniq*, and *vi\**.

10195 The *iconv()*, *iconv\_close()*, and *iconv\_open()* functions are available to allow an application to  
10196 convert character data between supported character sets.

10197 The *gencat* utility and the *catopen()*, *catclose()*, and *catgets()* functions for message catalog  
10198 manipulation are available on XSI-conformant systems.

**10199 Unsatisfied Requirements**

10200 Some aspects of multi-byte character and state-encoded character encodings have not yet been  
10201 addressed. The C-language functions, such as *getopt()*, are generally limited to single-byte  
10202 characters. The effect of the *LC\_MESSAGES* variable on message formats is only suggested at  
10203 this time.

**10204 D.2.16 Inter-User Communication**

10205 The following utilities address user requirements in this area: *cksum*, *mailx\**, *mesg\**, *patch\**, *pax*,  
10206 *talk\**, *uudecode\**, *uuencode\**, *who\**, and *write\**.

10207 The historical UUCP utilities are included on XSI-conformant systems.

**10208 Unsatisfied Requirements**

10209 None.

**10210 D.2.17 System Environment**

10211 The following utilities address user requirements in this area: *chgrp*, *chmod*, *chown*, *df\**, *du\**, *env*,  
10212 *getconf*, *id*, *logger*, *logname*, *mesg\**, *newgrp\**, *ps\**, *stty*, *tput\**, *tty*, *umask*, *uname*, and *who\**.

10213 The *closelog()*, *openlog()*, *setlogmask()*, and *syslog()* functions provide System Logging facilities  
10214 on XSI-conformant systems; these are analogous to the *logger* utility.

**10215 Unsatisfied Requirements**

10216 None.

**10217 D.2.18 Printing**

10218 The following utilities address user requirements in this area: *pr* and *lp*.

10219 **Unsatisfied Requirements**

10220 There are no features to control the formatting or scheduling of the print jobs.

10221 **D.2.19 Software Development**10222 The following utilities address user requirements in this area: *ar*, *asa*, *awk*, *c99*, *ctags\**, *fort77*,  
10223 *getconf*, *getopts*, *lex*, *localedef*, *make*, *nm\**, *od*, *patch\**, *pax*, *strings\**, *strip*, *time\**, and *yacc*.10224 The *system()*, *popen()*, *pclose()*, *regcomp()*, *regexec()*, *regerror()*, *regfree()*, *fnmatch()*, *getopt()*,  
10225 *glob()*, *globfree()*, *wordexp()*, and *wordfree()* functions allow C-language programmers to access  
10226 some of the interfaces used by the utilities, such as argument processing, regular expressions,  
10227 and pattern matching.10228 The SCCS source-code control system utilities are available on systems supporting the XSI  
10229 Development option.10230 **Unsatisfied Requirements**10231 There are no language-specific development tools related to languages other than C and  
10232 FORTRAN. The C tools are more complete and varied than the FORTRAN tools. There is no  
10233 data dictionary or other CASE-like development tools.10234 **D.2.20 Future Growth**10235 It is arguable whether or not all functionality to support applications is potentially within the  
10236 scope of IEEE Std. 1003.1-200x. As a simple matter of practicality, it cannot be. Areas such as  
10237 general networking, graphics, application domain-specific functionality, windowing, and so on,  
10238 should be in unique standards. As such, they are properly “Unsatisfied Requirements” in terms  
10239 of providing fully portable applications, but ones which are outside the scope of  
10240 IEEE Std. 1003.1-200x.10241 **D.3 Profiling Considerations**10242 This section offers guidance to writers of profiles on how the configurable options, limits, and  
10243 optional behavior of IEEE Std. 1003.1-200x should be cited in profiles. Profile writers should  
10244 consult the general guidance in POSIX.0 when writing POSIX Standardized Profiles.10245 The information in this section is an inclusive list of features that should be considered by profile  
10246 writers. Further subsetting of IEEE Std. 1003.1-200x, including the specification of behavior  
10247 currently described as unspecified, undefined, implementation-defined, or with the verbs “may”  
10248 or “need not” violates the intent of the developers of IEEE Std. 1003.1-200x and the guidelines of  
10249 ISO/IEC TR 10000-1. A set of profiling option groups is described in the Base Definitions  
10250 volume of IEEE Std. 1003.1-200x, based on the IEEE Std. 1003.13-1998 options, with the addition  
10251 of a new profiling option group called `_POSIX_NETWORKING`.

10252 **D.3.1 Configuration Options**

10253 There are two set of options suggested by IEEE Std. 1003.1-200x: those for POSIX-conforming  
 10254 systems and those for X/Open System Interface (XSI) conformance. The requirements for XSI  
 10255 conformance are documented in the Base Definitions volume of IEEE Std. 1003.1-200x and not  
 10256 discussed further here, as they superset the POSIX conformance requirements.

10257 **D.3.2 Configuration Options (Shell and Utilities)**

10258 There are three broad optional configurations for the Shell and Utilities volume of  
 10259 IEEE Std. 1003.1-200x: basic execution system, development system, and user portability  
 10260 interactive system. The options to support these, and other minor configuration options, are  
 10261 listed in the Base Definitions volume of IEEE Std. 1003.1-200x, Chapter 2, Conformance. Profile  
 10262 writers should consult the following list and the comments concerning user requirements  
 10263 addressed by various components in Section D.2 (on page 3560).

## 10264 POSIX2\_UPE

10265 The system supports the User Portability Utilities option.

10266 This option is a requirement for a user portability interactive system. It is required  
 10267 frequently except for those systems, such as embedded realtime or dedicated application  
 10268 systems, that support little or no interactive time-sharing work by users or operators. XSI-  
 10269 conformant systems support this option.

## 10270 POSIX2\_SW\_DEV

10271 The system supports the Software Development Utilities option.

10272 This option is required by many systems, even those in which actual software development  
 10273 does not occur. The *make* utility, in particular, is required by many application software  
 10274 packages as they are installed onto the system. If POSIX2\_C\_DEV is supported,  
 10275 POSIX2\_SW\_DEV is almost a mandatory requirement because of *ar* and *make*.

## 10276 POSIX2\_C\_BIND

10277 The system supports the C-Language Bindings option.

10278 This option is required on some systems developing complex C applications or on any  
 10279 system installing C applications in source form that require the functions in this option. The  
 10280 *system()* and *popen()* functions, in particular, are widely used by applications; the others are  
 10281 rather more specialized.

## 10282 POSIX2\_C\_DEV

10283 The system supports the C-Language Development Utilities option.

10284 This option is required by many systems, even those in which actual C-language software  
 10285 development does not occur. The *c99* utility, in particular, is required by many application  
 10286 software packages as they are installed onto the system. The *lex* and *yacc* utilities are used  
 10287 less frequently.

## 10288 POSIX2\_FORT\_DEV

10289 The system supports the FORTRAN Development Utilities option

10290 As with C, this option is needed on any system developing or installing FORTRAN  
 10291 applications in source form.

## 10292 POSIX2\_FORT\_RUN

10293 The system supports the FORTRAN Runtime Utilities option.

10294 This option is required for some FORTRAN applications that need the *asa* utility to convert  
 10295 Hollerith printing statement output. It is unknown how frequently this occurs.

- 10296 POSIX2\_LOCALEDEF  
 10297 The system supports the creation of locales.
- 10298 This option is needed if applications require their own customized locale definitions to  
 10299 operate. It is presently unknown whether many applications are dependent on this.  
 10300 However, the option is virtually mandatory for systems in which internationalized  
 10301 applications are developed.
- 10302 XSI-conformant systems support this option.
- 10303 POSIX2\_PBS  
 10304 The system supports the Batch Environment option.
- 10305 POSIX2\_PBS\_ACCOUNTING  
 10306 The system supports the optional feature of accounting within the Batch Environment  
 10307 option. It will be required in servers that implement the optional feature of accounting.
- 10308 POSIX2\_PBS\_CHECKPOINT  
 10309 The systems supports the optional feature of checkpoint/restart within the Batch  
 10310 Environment option.
- 10311 POSIX2\_PBS\_LOCATE  
 10312 The system supports the optional feature of locating batch jobs within the Batch  
 10313 Environment option.
- 10314 POSIX2\_PBS\_MESSAGE  
 10315 The system supports the optional feature of sending messages to batch jobs within the  
 10316 Batch Environment option.
- 10317 POSIX2\_PBS\_TRACK  
 10318 The system supports the optional feature of tracking batch jobs within the Batch  
 10319 Environment option.
- 10320 POSIX2\_CHAR\_TERM  
 10321 The system supports at least one terminal type capable of all operations described in  
 10322 IEEE Std. 1003.1-200x.
- 10323 On systems with POSIX2\_UPE, this option is almost always required. It was developed  
 10324 solely to allow certain specialized vendors and user applications to bypass the requirement  
 10325 for general-purpose asynchronous terminal support. For example, an application and  
 10326 system that was suitable for block-mode terminals, such as IBM 3270s, would not need this  
 10327 option.
- 10328 XSI-conformant systems support this option.

### 10329 D.3.3 Configurable Limits

- 10330 Very few of the limits need to be increased for profiles. No profile can cite lower values.
- 10331 {POSIX2\_BC\_BASE\_MAX}  
 10332 {POSIX2\_BC\_DIM\_MAX}  
 10333 {POSIX2\_BC\_SCALE\_MAX}  
 10334 {POSIX2\_BC\_STRING\_MAX}  
 10335 No increase is anticipated for any of these *bc* values, except for very specialized applications  
 10336 involving huge numbers.
- 10337 {POSIX2\_COLL\_WEIGHTS\_MAX}  
 10338 Some natural languages with complex collation requirements require an increase from the  
 10339 default 2 to 4; no higher numbers are anticipated.

10340 {POSIX2\_EXPR\_NEST\_MAX}  
 10341 No increase is anticipated.

10342 {POSIX2\_LINE\_MAX}  
 10343 This number is much larger than most historical applications have been able to use. At some  
 10344 future time, applications may be rewritten to take advantage of even larger values.

10345 {POSIX2\_RE\_DUP\_MAX}  
 10346 No increase is anticipated.

10347 {POSIX2\_VERSION}  
 10348 This is actually not a limit, but a standard version stamp. Generally, a profile should specify  
 10349 Shell and Utilities volume of IEEE Std. 1003.1-200x, Chapter 2, Shell Command Language by  
 10350 name in the normative references section, not this value.

#### 10351 **D.3.4 Configuration Options (System Interfaces)**

10352 {NGROUPS\_MAX}  
 10353 A non-zero value indicates that the implementation supports supplementary groups.

10354 This option is needed where there is a large amount of shared use of files, but where a  
 10355 certain amount of protection is needed. Many profiles<sup>3</sup> are known to require this option; it  
 10356 should only be required if needed, but it should never be prohibited.

10357 \_POSIX\_ADVISORY\_INFO  
 10358 The system provides advisory information for file management.

10359 This option allows the application to specify advisory information that can be used to  
 10360 achieve better or even deterministic response time in file manager or input and output  
 10361 operations.

10362 \_POSIX\_ASYNCHRONOUS\_IO  
 10363 The system provides concurrent process execution and input and output transfers.

10364 This option was created to support historical systems that did not provide the feature. It  
 10365 should only be required if needed, but it should never be prohibited.

10366 \_POSIX\_BARRIERS  
 10367 The system supports barrier synchronization.

10368 This option was created to allow efficient synchronization of multiple parallel threads in  
 10369 multi-processor systems in which the operation is supported in part by the hardware  
 10370 architecture.

10371 \_POSIX\_CHOWN\_RESTRICTED  
 10372 The system restricts the right to “give away” files to other users.

10373 This option should be carefully investigated before it is required. Some applications expect  
 10374 that they can change the ownership of files in this way. It is provided where either security  
 10375 or system account requirements cause this ability to be a problem. It is also known to be  
 10376 specified in many profiles.

10377 \_\_\_\_\_  
 10378 3. There are no formally approved profiles of IEEE Std. 1003.1-200x at the time of publication; the reference here is to various  
 10379 profiles generated by private bodies or governments.

- 10380     \_POSIX\_CLOCK\_SELECTION  
10381         The system supports the Clock Selection option.
- 10382         This option allows applications to request a high resolution sleep in order to suspend a  
10383         thread during a relative time interval, or until an absolute time value, using the desired  
10384         clock. It also allows the application to select the clock used in a *pthread\_cond\_timedwait()*  
10385         function call.
- 10386     \_POSIX\_CPUTIME  
10387         The system supports the Process CPU-Time Clocks option.
- 10388         This option allows applications to use a new clock that measures the execution times of  
10389         processes or threads, and the possibility to create timers based upon these clocks, for  
10390         runtime detection (and treatment) of execution time overruns.
- 10391     \_POSIX\_FSYNC  
10392         The system supports file synchronization requests.
- 10393         This option was created to support historical systems that did not provide the feature.  
10394         Applications that are expecting guaranteed completion of their input and output operations  
10395         should require the *\_POSIX\_SYNC\_IO* option. This option should never be prohibited.
- 10396         XSI-conformant systems support this option.
- 10397     \_POSIX\_IPV6  
10398         The system supports facilities related to Internet Protocol Version 6 (IPv6).
- 10399         This option was created to allow systems to transition to IPv6.
- 10400     \_POSIX\_JOB\_CONTROL  
10401         Job control facilities are mandatory in IEEE Std. 1003.1-200x.
- 10402         The option was created primarily to support historical systems that did not provide the  
10403         feature. Many existing profiles now require it; it should only be required if needed, but it  
10404         should never be prohibited. Most applications that use it can run when it is not present,  
10405         although with a degraded level of user convenience.
- 10406     \_POSIX\_MAPPED\_FILES  
10407         The system supports the mapping of regular files into the process address space.
- 10408         XSI-conformant systems support this option.
- 10409         Both this option and the Shared Memory Objects option provide shared access to memory  
10410         objects in the process address space. The functions defined under this option provide the  
10411         functionality of existing practice for mapping regular files. This functionality was deemed  
10412         unnecessary, if not inappropriate, for embedded systems applications and, hence, is  
10413         provided under this option. It should only be required if needed, but it should never be  
10414         prohibited.
- 10415     \_POSIX\_MEMLOCK  
10416         The system supports the locking of the address space.
- 10417         This option was created to support historical systems that did not provide the feature. It  
10418         should only be required if needed, but it should never be prohibited.
- 10419     \_POSIX\_MEMLOCK\_RANGE  
10420         The system supports the locking of specific ranges of the address space.
- 10421         For applications that have well-defined sections that need to be locked and others that do  
10422         not, IEEE Std. 1003.1-200x supports an optional set of functions to lock or unlock a range of  
10423         process addresses. The following are two reasons for having a means to lock down a

- 10424 specific range:
- 10425 1. An asynchronous event handler function that must respond to external events in a  
10426 deterministic manner such that page faults cannot be tolerated
- 10427 2. An input/output “buffer” area that is the target for direct-to-process I/O, and the  
10428 overhead of implicit locking and unlocking for each I/O call cannot be tolerated
- 10429 It should only be required if needed, but it should never be prohibited.
- 10430 \_POSIX\_MEMORY\_PROTECTION  
10431 The system supports memory protection.
- 10432 XSI-conformant systems support this option.
- 10433 The provision of this option typically imposes additional hardware requirements. It should  
10434 never be prohibited.
- 10435 \_POSIX\_PRIORITIZED\_IO  
10436 The system provides prioritization for input and output operations.
- 10437 The use of this option may interfere with the ability of the system to optimize input and  
10438 output throughput. It should only be required if needed, but it should never be prohibited.
- 10439 \_POSIX\_MESSAGE\_PASSING  
10440 The system supports the passing of messages between processes.
- 10441 This option was created to support historical systems that did not provide the feature. The  
10442 functionality adds a high-performance XSI interprocess communication facility for local  
10443 communication. It should only be required if needed, but it should never be prohibited.
- 10444 \_POSIX\_MONOTONIC\_CLOCK  
10445 The system supports the Monotonic Clock option.
- 10446 This option allows realtime applications to rely on a monotonically increasing clock that  
10447 does not jump backwards, and whose value does not change except for the regular ticking  
10448 of the clock.
- 10449 \_POSIX\_PRIORITY\_SCHEDULING  
10450 The system provides priority-based process scheduling.
- 10451 Support of this option provides predictable scheduling behavior, allowing applications to  
10452 determine the order in which processes that are ready to run are granted access to a  
10453 processor. It should only be required if needed, but it should never be prohibited.
- 10454 \_POSIX\_REALTIME\_SIGNALS  
10455 The system provides prioritized, queued signals with associated data values.
- 10456 This option was created to support historical systems that did not provide the features. It  
10457 should only be required if needed, but it should never be prohibited.
- 10458 \_POSIX\_REGEX  
10459 Support for regular expression facilities are mandatory in IEEE Std. 1003.1-200x.
- 10460 \_POSIX\_SAVED\_IDS  
10461 Support for this feature is mandatory in IEEE Std. 1003.1-200x.
- 10462 Certain classes of applications rely on it for proper operation, and there is no alternative  
10463 short of giving the application root privileges on most implementations that did not provide  
10464 \_POSIX\_SAVED\_IDS.



- 10465 \_POSIX\_SEMAPHORES  
10466 The system provides counting semaphores.
- 10467 This option was created to support historical systems that did not provide the feature. It  
10468 should only be required if needed, but it should never be prohibited.
- 10469 \_POSIX\_SHARED\_MEMORY\_OBJECTS  
10470 The system supports the mapping of shared memory objects into the process address space.
- 10471 Both this option and the Memory Mapped Files option provide shared access to memory  
10472 objects in the process address space. The functions defined under this option provide the  
10473 functionality of existing practice for shared memory objects. This functionality was deemed  
10474 appropriate for embedded systems applications and, hence, is provided under this option. It  
10475 should only be required if needed, but it should never be prohibited.
- 10476 \_POSIX\_SHELL  
10477 Support for the *sh* utility command line interpreter is mandatory in IEEE Std. 1003.1-200x.
- 10478 \_POSIX\_SPAWN  
10479 The system supports the spawn option.
- 10480 This option provides applications with an efficient mechanism to spawn execution of a new  
10481 process.
- 10482 \_POSIX\_SPINLOCKS  
10483 The system supports spin locks.
- 10484 This option was created to support a simple and efficient synchronization mechanism for  
10485 threads executing in multi-processor systems.
- 10486 \_POSIX\_SPORADIC\_SERVER  
10487 The system supports the sporadic server scheduling policy.
- 10488 This option provides applications with a new scheduling policy for scheduling aperiodic  
10489 processes or threads in hard realtime applications.
- 10490 \_POSIX\_SYNCHRONIZED\_IO  
10491 The system supports guaranteed file synchronization.
- 10492 This option was created to support historical systems that did not provide the feature.  
10493 Applications that are expecting guaranteed completion of their input and output operations  
10494 should require this option, rather than the File Synchronization option. It should only be  
10495 required if needed, but it should never be prohibited.
- 10496 \_POSIX\_THREADS  
10497 The system supports multiple threads of control within a single process.
- 10498 This option was created to support historical systems that did not provide the feature.  
10499 Applications written assuming a multi-threaded environment would be expected to require  
10500 this option. It should only be required if needed, but it should never be prohibited.
- 10501 XSI-conformant systems support this option.
- 10502 \_POSIX\_THREAD\_ATTR\_STACKADDR  
10503 The system supports specification of the stack address for a created thread.
- 10504 Applications may take advantage of support of this option for performance benefits, but  
10505 dependence on this feature should be minimized. This option should never be prohibited.
- 10506 XSI-conformant systems support this option.

- 10507 \_POSIX\_THREAD\_ATTR\_STACKSIZE  
10508 The system supports specification of the stack size for a created thread.
- 10509 Applications may require this option in order to ensure proper execution, but such usage  
10510 limits portability and dependence on this feature should be minimized. It should only be  
10511 required if needed, but it should never be prohibited.
- 10512 XSI-conformant systems support this option.
- 10513 \_POSIX\_THREAD\_PRIORITY\_SCHEDULING  
10514 The system provides priority-based thread scheduling.
- 10515 Support of this option provides predictable scheduling behavior, allowing applications to  
10516 determine the order in which threads that are ready to run are granted access to a processor.  
10517 It should only be required if needed, but it should never be prohibited.
- 10518 \_POSIX\_THREAD\_PRIO\_INHERIT  
10519 The system provides mutual exclusion operations with priority inheritance.
- 10520 Support of this option provides predictable scheduling behavior, allowing applications to  
10521 determine the order in which threads that are ready to run are granted access to a processor.  
10522 It should only be required if needed, but it should never be prohibited.
- 10523 \_POSIX\_THREAD\_PRIO\_PROTECT  
10524 The system supports a priority ceiling emulation protocol for mutual exclusion operations.
- 10525 Support of this option provides predictable scheduling behavior, allowing applications to  
10526 determine the order in which threads that are ready to run are granted access to a processor.  
10527 It should only be required if needed, but it should never be prohibited.
- 10528 \_POSIX\_THREAD\_PROCESS\_SHARED  
10529 The system provides shared access among multiple processes to synchronization objects.
- 10530 This option was created to support historical systems that did not provide the feature. It  
10531 should only be required if needed, but it should never be prohibited.
- 10532 XSI-conformant systems support this option.
- 10533 \_POSIX\_THREAD\_SAFE\_FUNCTIONS  
10534 The system provides thread-safe versions of all of the POSIX.1 functions.
- 10535 This option is required if the Threads option is supported. This is a separate option because  
10536 thread-safe functions are useful in implementations providing other mechanisms for  
10537 concurrency. It should only be required if needed, but it should never be prohibited.
- 10538 XSI-conformant systems support this option.
- 10539 \_POSIX\_THREAD\_SPORADIC\_SERVER  
10540 The system supports the thread sporadic server scheduling policy.
- 10541 Support for this option provides applications with a new scheduling policy for scheduling  
10542 aperiodic threads in hard realtime applications.
- 10543 \_POSIX\_TIMEOUTS  
10544 The system provides timeouts for some blocking services.
- 10545 This option was created to provide a timeout capability to system services, thus allowing  
10546 applications to include better error detection, and recovery capabilities.
- 10547 \_POSIX\_TIMERS  
10548 The system provides higher resolution clocks with multiple timers per process.

10549 This option was created to support historical systems that did not provide the features. This  
 10550 option is appropriate for applications requiring higher resolution timestamps or needing to  
 10551 control the timing of multiple activities. It should only be required if needed, but it should  
 10552 never be prohibited.

10553 `_POSIX_TRACE`

10554 The system supports the trace option.

10555 This option was created to allow applications to perform tracing.

10556 `_POSIX_TRACE_EVENT_FILTER`

10557 The system supports the trace event filter option.

10558 This option is dependent on support of the Trace option.

10559 `_POSIX_TRACE_INHERIT`

10560 The system supports the trace inherit option.

10561 This option is dependent on support of the Trace option.

10562 `_POSIX_TRACE_LOG`

10563 The system supports the trace log option.

10564 This option is dependent on support of the Trace option.

10565 `_POSIX_TYPED_MEMORY_OBJECTS`

10566 The system supports typed memory objects.

10567 This option was created to allow realtime applications to access different kinds of physical  
 10568 memory, and allow processes in these applications to share portions of this memory.

### 10569 D.3.5 Configurable Limits

10570 In general, the configurable limits in the `<limits.h>` header defined in the Base Definitions  
 10571 volume of IEEE Std. 1003.1-200x have been set to minimal values; many applications or  
 10572 implementations may require larger values. No profile can cite lower values.

10573 `{AIO_LISTIO_MAX}`

10574 The current minimum is likely to be inadequate for most applications. It is expected that  
 10575 this value will be increased by profiles requiring support for list input and output  
 10576 operations.

10577 `{AIO_MAX}`

10578 The current minimum is likely to be inadequate for most applications. It is expected that  
 10579 this value will be increased by profiles requiring support for asynchronous input and  
 10580 output operations.

10581 `{AIO_PRIO_DELTA_MAX}`

10582 The functionality associated with this limit is needed only by sophisticated applications. It  
 10583 is not expected that this limit would need to be increased under a general-purpose profile.

10584 `{ARG_MAX}`

10585 The current minimum is likely to need to be increased for profiles, particularly as larger  
 10586 amounts of information are passed through the environment. Many implementations are  
 10587 believed to support larger values.

10588 `{CHILD_MAX}`

10589 The current minimum is suitable only for systems where a single user is not running  
 10590 applications in parallel. It is significantly too low for any system also requiring windows,  
 10591 and if `_POSIX_JOB_CONTROL` is specified, it should be raised.

- 10592 {CLOCKRES\_MIN}  
10593 It is expected that profiles will require a finer granularity clock, perhaps as fine as 1  $\mu$ s,  
10594 represented by a value of 1 000 for this limit.
- 10595 {DELAYTIMER\_MAX}  
10596 It is believed that most implementations will provide larger values.
- 10597 {LINK\_MAX}  
10598 For most applications and usage, the current minimum is adequate. Many implementations  
10599 have a much larger value, but this should not be used as a basis for raising the value unless  
10600 the applications to be used require it.
- 10601 {LOGIN\_NAME\_MAX}  
10602 This is not actually a limit, but an implementation parameter. No profile should impose a  
10603 requirement on this value.
- 10604 {MAX\_CANON}  
10605 For most purposes, the current minimum is adequate. Unless high-speed burst serial  
10606 devices are used, it should be left as is.
- 10607 {MAX\_INPUT}  
10608 See {MAX\_CANON}.
- 10609 {MQ\_OPEN\_MAX}  
10610 The current minimum should be adequate for most profiles.
- 10611 {MQ\_PRIO\_MAX}  
10612 The current minimum corresponds to the required number of process scheduling priorities.  
10613 Many realtime practitioners believe that the number of message priority levels ought to be  
10614 the same as the number of execution scheduling priorities.
- 10615 {NAME\_MAX}  
10616 Many implementations now support larger values, and many applications and users  
10617 assume that larger names can be used. Many existing profiles also specify a larger value.  
10618 Specifying this value will reduce the number of conforming implementations, although this  
10619 might not be a significant consideration over time. Values greater than 255 should not be  
10620 required.
- 10621 {NGROUPS\_MAX}  
10622 The value selected will typically be 8 or larger.
- 10623 {OPEN\_MAX}  
10624 The historically common value for this has been 20. Many implementations support larger  
10625 values. If applications that use larger values are anticipated, an appropriate value should be  
10626 specified.
- 10627 {PAGESIZE}  
10628 This is not actually a limit, but an implementation parameter. No profile should impose a  
10629 requirement on this value.
- 10630 {PATH\_MAX}  
10631 Historically, the minimum has been either 1024 or indefinite, depending on the  
10632 implementation. Few applications actually require values larger than 256, but some users  
10633 may create file hierarchies that must be accessed with longer paths. This value should only  
10634 be changed if there is a clear requirement.
- 10635 {PIPE\_BUF}  
10636 The current minimum is adequate for most applications. Historically, it has been larger. If  
10637 applications that write single transactions larger than this are anticipated, it should be

10638           increased. Applications that write lines of text larger than this probably do not need it  
10639           increased, as the text line is delimited by a newline.

10640           {POSIX\_VERSION}  
10641           This is actually not a limit, but a standard version stamp. Generally, a profile should specify  
10642           IEEE Std. 1003.1-200x by a name in the normative references section, not this value.

10643           {PTHREAD\_DESTRUCTOR\_ITERATIONS}  
10644           It is unlikely that applications will need larger values to avoid loss of memory resources.

10645           {PTHREAD\_KEYS\_MAX}  
10646           The current value should be adequate for most profiles.

10647           {PTHREAD\_STACK\_MIN}  
10648           This should not be treated as an actual limit, but as an implementation parameter. No  
10649           profile should impose a requirement on this value.

10650           {PTHREAD\_THREADS\_MAX}  
10651           It is believed that most implementations will provide larger values.

10652           {RTSIG\_MAX}  
10653           The current limit was chosen so that the set of POSIX.1 signal numbers can fit within a 32-  
10654           bit field. It is recognized that most existing implementations define many more signals than  
10655           are specified in POSIX.1 and, in fact, many implementations have already exceeded 32  
10656           signals (including the “null signal”). Support of {\_POSIX\_RTSIG\_MAX} additional signals  
10657           may push some implementations over the single 32-bit word line, but is unlikely to push  
10658           any implementations that are already over that line beyond the 64 signal line.

10659           {SEM\_NSEMS\_MAX}  
10660           The current value should be adequate for most profiles.

10661           {SEM\_VALUE\_MAX}  
10662           The current value should be adequate for most profiles.

10663           {SSIZE\_MAX}  
10664           This limit reflects fundamental hardware characteristics (the size of an integer), and should  
10665           not be specified unless it is clearly required. Extreme care should be taken to assure that  
10666           any value that might be specified does not unnecessarily eliminate implementations  
10667           because of accidents of hardware design.

10668           {STREAM\_MAX}  
10669           This limit is very closely related to {OPEN\_MAX}. It should never be larger than  
10670           {OPEN\_MAX}, but could reasonably be smaller for application areas where most files are  
10671           not accessed through *stdio*. Some implementations may limit {STREAM\_MAX} to 20 but  
10672           allow {OPEN\_MAX} to be considerably larger. Such implementations should be allowed for  
10673           if the applications permit.

10674           {TIMER\_MAX}  
10675           The current limit should be adequate for most profiles, but it may need to be larger for  
10676           applications with a large number of asynchronous operations.

10677           {TTY\_NAME\_MAX}  
10678           This is not actually a limit, but an implementation parameter. No profile should impose a  
10679           requirement on this value.

10680           {TZNAME\_MAX}  
10681           The minimum has been historically adequate, but if longer timezone names are anticipated  
10682           (particularly such values as UTC-1), this should be increased.

**10683 D.3.6 Optional Behavior**

10684 In IEEE Std. 1003.1-200x, there are no instances of the terms unspecified, undefined,  
10685 implementation-defined, or with the verbs “may” or “need not”, that the developers of  
10686 IEEE Std. 1003.1-200x anticipate or sanction as suitable for profile or test method citation. All of  
10687 these are merely warnings to portable applications to avoid certain areas that can vary from  
10688 system to system, and even over time on the same system. In many cases, these terms are used  
10689 explicitly to support extensions, but profiles should not anticipate and require such extensions;  
10690 future versions of IEEE Std. 1003.1-200x may do so.