



Multithreaded Programming Guide

Sun Microsystems, Inc.
4150 Network Circle
Santa Clara, CA 95054
U.S.A.

Part No: 806-6867-10
May 2002

Copyright 2002 Sun Microsystems, Inc. 4150 Network Circle, Santa Clara, CA 95054 U.S.A. All rights reserved.

This product or document is protected by copyright and distributed under licenses restricting its use, copying, distribution, and decompilation. No part of this product or document may be reproduced in any form by any means without prior written authorization of Sun and its licensors, if any. Third-party software, including font technology, is copyrighted and licensed from Sun suppliers.

Parts of the product may be derived from Berkeley BSD systems, licensed from the University of California. UNIX is a registered trademark in the U.S. and other countries, exclusively licensed through X/Open Company, Ltd.

Sun, Sun Microsystems, the Sun logo, docs.sun.com, AnswerBook, AnswerBook2, and Solaris are trademarks, registered trademarks, or service marks of Sun Microsystems, Inc. in the U.S. and other countries. All SPARC trademarks are used under license and are trademarks or registered trademarks of SPARC International, Inc. in the U.S. and other countries. Products bearing SPARC trademarks are based upon an architecture developed by Sun Microsystems, Inc.

The OPEN LOOK and Sun™ Graphical User Interface was developed by Sun Microsystems, Inc. for its users and licensees. Sun acknowledges the pioneering efforts of Xerox in researching and developing the concept of visual or graphical user interfaces for the computer industry. Sun holds a non-exclusive license from Xerox to the Xerox Graphical User Interface, which license also covers Sun's licensees who implement OPEN LOOK GUIs and otherwise comply with Sun's written license agreements.

Federal Acquisitions: Commercial Software—Government Users Subject to Standard License Terms and Conditions.

DOCUMENTATION IS PROVIDED "AS IS" AND ALL EXPRESS OR IMPLIED CONDITIONS, REPRESENTATIONS AND WARRANTIES, INCLUDING ANY IMPLIED WARRANTY OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE OR NON-INFRINGEMENT, ARE DISCLAIMED, EXCEPT TO THE EXTENT THAT SUCH DISCLAIMERS ARE HELD TO BE LEGALLY INVALID.

Copyright 2002 Sun Microsystems, Inc. 4150 Network Circle, Santa Clara, CA 95054 U.S.A. Tous droits réservés

Ce produit ou document est protégé par un copyright et distribué avec des licences qui en restreignent l'utilisation, la copie, la distribution, et la décompilation. Aucune partie de ce produit ou document ne peut être reproduite sous aucune forme, par quelque moyen que ce soit, sans l'autorisation préalable et écrite de Sun et de ses bailleurs de licence, s'il y en a. Le logiciel détenu par des tiers, et qui comprend la technologie relative aux polices de caractères, est protégé par un copyright et licencié par des fournisseurs de Sun.

Des parties de ce produit pourront être dérivées du système Berkeley BSD licenciés par l'Université de Californie. UNIX est une marque déposée aux Etats-Unis et dans d'autres pays et licenciée exclusivement par X/Open Company, Ltd.

Sun, Sun Microsystems, le logo Sun, docs.sun.com, AnswerBook, AnswerBook2, et Solaris sont des marques de fabrique ou des marques déposées, ou marques de service, de Sun Microsystems, Inc. aux Etats-Unis et dans d'autres pays. Toutes les marques SPARC sont utilisées sous licence et sont des marques de fabrique ou des marques déposées de SPARC International, Inc. aux Etats-Unis et dans d'autres pays. Les produits portant les marques SPARC sont basés sur une architecture développée par Sun Microsystems, Inc.

L'interface d'utilisation graphique OPEN LOOK et Sun™ a été développée par Sun Microsystems, Inc. pour ses utilisateurs et licenciés. Sun reconnaît les efforts de pionniers de Xerox pour la recherche et le développement du concept des interfaces d'utilisation visuelle ou graphique pour l'industrie de l'informatique. Sun détient une licence non exclusive de Xerox sur l'interface d'utilisation graphique Xerox, cette licence couvrant également les licenciés de Sun qui mettent en place l'interface d'utilisation graphique OPEN LOOK et qui en outre se conforment aux licences écrites de Sun.

CETTE PUBLICATION EST FOURNIE "EN L'ETAT" ET AUCUNE GARANTIE, EXPRESSE OU IMPLICITE, N'EST ACCORDEE, Y COMPRIS DES GARANTIES CONCERNANT LA VALEUR MARCHANDE, L'APTITUDE DE LA PUBLICATION A REpondre A UNE UTILISATION PARTICULIERE, OU LE FAIT QU'ELLE NE SOIT PAS CONTREFAISANTE DE PRODUIT DE TIERS. CE DENI DE GARANTIE NE S'APPLIQUERAIT PAS, DANS LA MESURE OU IL SERAIT TENU JURIDIQUEMENT NUL ET NON AVENU.



020115@3062



Contents

Preface 11

1	Covering Multithreading Basics	15
	Defining Multithreading Terms	15
	Meeting Multithreading Standards	17
	Benefiting From Multithreading	17
	Improving Application Responsiveness	17
	Using Multiprocessors Efficiently	17
	Improving Program Structure	18
	Using Fewer System Resources	18
	Combining Threads and RPC	18
	Understanding Basic Multithreading Concepts	18
	Concurrency and Parallelism	18
	Looking at Multithreading Structure	19
	Scheduling	21
	Cancellation	22
	Synchronization	22
	Using the 64-bit Architecture	23
2	Basic Threads Programming	25
	The Threads Library	25
	Create a Default Thread	26
	Wait for Thread Termination	27
	A Simple Threads Example	28
	Detaching a Thread	29

Create a Key for Thread-Specific Data	30
Delete the Thread-Specific Data Key	31
Set Thread-Specific Data	32
Get Thread-Specific Data	33
Get the Thread Identifier	36
Compare Thread IDs	36
Initializing Threads	37
Yield Thread Execution	37
Set the Thread Priority	38
Get the Thread Priority	39
Send a Signal to a Thread	39
Access the Signal Mask of the Calling Thread	40
Forking Safely	41
Terminate a Thread	41
Finishing Up	42
Cancellation	43
Cancel a Thread	44
Enable or Disable Cancellation	45
Set Cancellation Type	45
Create a Cancellation Point	46
Push a Handler Onto the Stack	47
Pull a Handler Off the Stack	47
3 Thread Create Attributes	49
Attributes	50
Initialize Attributes	50
Destroy Attributes	52
Set Detach State	52
Get Detach State	53
Set Stack Guard Size	54
Get Stack Guard Size	55
Set Scope	56
Get Scope	57
Set Thread Concurrency Level	57
Get Thread Concurrency Level	58
Set Scheduling Policy	58
Get Scheduling Policy	59

	Set Inherited Scheduling Policy	60
	Get Inherited Scheduling Policy	61
	Set Scheduling Parameters	61
	Get Scheduling Parameters	62
	Set Stack Size	64
	Get Stack Size	65
	About Stacks	65
	Set Stack Address	67
	Get Stack Address	68
4	Programming with Synchronization Objects	69
	Mutual Exclusion Lock Attributes	70
	Initialize a Mutex Attribute Object	72
	Destroy a Mutex Attribute Object	72
	Set the Scope of a Mutex	73
	Get the Scope of a Mutex	74
	Set the Mutex Type Attribute	75
	Get the Mutex Type Attribute	76
	Set Mutex Attribute's Protocol	76
	Get Mutex Attribute's Protocol	79
	Set Mutex Attribute's Priority Ceiling	79
	Get Mutex Attribute's Priority Ceiling	80
	Set Mutex's Priority Ceiling	82
	Get Mutex's Priority Ceiling	83
	Set Mutex's Robust Attribute	83
	Get Mutex's Robust Attribute	85
	Using Mutual Exclusion Locks	86
	Initialize a Mutex	87
	Make Mutex Consistent	88
	Lock a Mutex	89
	Unlock a Mutex	91
	Lock With a Nonblocking Mutex	92
	Destroy a Mutex	93
	Mutex Lock Code Examples	94
	Condition Variable Attributes	98
	Initialize a Condition Variable Attribute	99
	Remove a Condition Variable Attribute	100

Set the Scope of a Condition Variable	101
Get the Scope of a Condition Variable	102
Using Condition Variables	102
Initialize a Condition Variable	103
Block on a Condition Variable	104
Unblock One Thread	105
Block Until a Specified Time	107
Block For a Specified Interval	108
Unblock All Threads	109
Destroy Condition Variable State	110
The Lost Wake-Up Problem	111
The Producer/Consumer Problem	111
Semaphores	114
Counting Semaphores	116
Initialize a Semaphore	116
Named Semaphores	118
Increment a Semaphore	118
Block on a Semaphore Count	119
Decrement a Semaphore Count	119
Destroy the Semaphore State	120
The Producer/Consumer Problem, Using Semaphores	121
Read-Write Lock Attributes	122
Initialize a Read-Write Lock Attribute	123
Destroy a Read-Write Lock Attribute	123
Set a Read-Write Lock Attribute	124
Get a Read-Write Lock Attribute	125
Using Read-Write Locks	125
Initialize a Read-Write Lock	126
Read Lock on Read-Write Lock	127
Read Lock With a Nonblocking Read-Write Lock	128
Write Lock on Read-Write Lock	128
Write Lock With a Nonblocking Read-Write Lock	129
Unlock a Read-Write Lock	129
Destroy a Read-Write Lock	130
Synchronization Across Process Boundaries	131
Producer/Consumer Problem Example	131
Interprocess Locking Without the Threads Library	133
Comparing Primitives	133

5	Programming With the Operating Environment	135
	Process Creation—Forking Issues	135
	The Fork-One Model	136
	The Fork-All Model	139
	Choosing the Right Fork	139
	Process Creation—exec(2) and exit(2) Issues	139
	Timers, Alarms, and Profiling	140
	Per-LWP POSIX Timers	140
	Per-Thread Alarms	141
	Profiling	141
	Nonlocal Goto—setjmp(3C) and longjmp(3C)	142
	Resource Limits	142
	LWPs and Scheduling Classes	142
	Timeshare Scheduling	143
	Realtime Scheduling	143
	Fair Share Scheduling	144
	Fixed Priority Scheduling	144
	Extending Traditional Signals	145
	Synchronous Signals	146
	Asynchronous Signals	146
	Continuation Semantics	146
	Operations on Signals	147
	Thread-Directed Signals	149
	Completion Semantics	150
	Signal Handlers and Async-Signal Safety	151
	Interrupted Waits on Condition Variables	153
	I/O Issues	154
	I/O as a Remote Procedure Call	154
	Tamed Asynchrony	155
	Asynchronous I/O	155
	Shared I/O and New I/O System Calls	156
	Alternatives to getc(3C) and putc(3C)	157
6	Safe and Unsafe Interfaces	159
	Thread Safety	159
	MT Interface Safety Levels	160
	Reentrant Functions for Unsafe Interfaces	162

Async-Signal-Safe Functions	162
MT Safety Levels for Libraries	163
Unsafe Libraries	164
7 Compiling and Debugging	165
Compiling a Multithreaded Application	165
Preparing for Compilation	165
Choosing Solaris or POSIX Semantics	166
Including <thread.h> or <pthread.h>	166
Defining _REENTRANT or _POSIX_C_SOURCE	167
Linking With libthread or libpthread	167
Linking With -lrt for POSIX Semaphores	168
Link Old With New	169
The Alternate libthread	169
Debugging a Multithreaded Program	169
Common Oversights	169
Tracing and Debugging With the TNF Utilities	170
Using truss(1)	171
Using mdb(1)	171
Using dbx	171
8 Programming With Solaris Threads	175
Comparing APIs for Solaris Threads and POSIX Threads	175
Major API Differences	176
Function Comparison Table	176
Unique Solaris Threads Functions	180
Suspend Thread Execution	180
Continue a Suspended Thread	181
Similar Synchronization Functions—Read-Write Locks	181
Initialize a Read-Write Lock	182
Acquire a Read Lock	183
Try to Acquire a Read Lock	184
Acquire a Write Lock	185
Try to Acquire a Write Lock	185
Unlock a Read-Write Lock	186
Destroy Read-Write Lock State	186
Similar Solaris Threads Functions	188

Create a Thread	188
Get the Minimal Stack Size	191
Get the Thread Identifier	191
Yield Thread Execution	192
Send a Signal to a Thread	192
Access the Signal Mask of the Calling Thread	192
Terminate a Thread	192
Wait for Thread Termination	193
Create a Thread-Specific Data Key	194
Set Thread-Specific Data	194
Get Thread-Specific Data	194
Set the Thread Priority	195
Get the Thread Priority	195
Similar Synchronization Functions—Mutual Exclusion Locks	196
Initialize a Mutex	196
Destroy a Mutex	197
Acquire a Mutex	198
Release a Mutex	198
Try to Acquire a Mutex	198
Similar Synchronization Functions—Condition Variables	199
Initialize a Condition Variable	199
Destroy a Condition Variable	200
Wait for a Condition	200
Wait for an Absolute Time	201
Wait for a Time Interval	201
Unblock One Thread	201
Unblock All Threads	202
Similar Synchronization Functions—Semaphores	202
Initialize a Semaphore	202
Increment a Semaphore	203
Block on a Semaphore Count	204
Decrement a Semaphore Count	204
Destroy the Semaphore State	204
Synchronization Across Process Boundaries	205
Using LWPs Between Processes	205
Producer/Consumer Problem Example	206
Special Issues for fork() and Solaris Threads	207

9	Programming Guidelines	209
	Rethinking Global Variables	209
	Providing for Static Local Variables	210
	Synchronizing Threads	211
	Single-Threaded Strategy	212
	Reentrance	212
	Avoiding Deadlock	214
	Deadlocks Related to Scheduling	215
	Locking Guidelines	215
	Following Some Basic Guidelines	216
	Creating and Using Threads	217
	Lightweight Processes	217
	Unbound Threads	219
	Bound Threads	219
	Thread Creation Guidelines	219
	Working With Multiprocessors	219
	The Underlying Architecture	220
	Summary	224
	Further Reading	224
A	Sample Application—Multithreaded grep	225
	Description of <code>tgrep</code>	225
	Getting Online Source Code	226
B	Solaris Threads Example: <code>barrier.c</code>	251
	Index	255

Preface

The *Multithreaded Programming Guide* describes the multithreaded programming interfaces for POSIX and Solaris threads in the Solaris™ Operating Environment. This guide shows application programmers how to create new multithreaded programs and how to add multithreading to existing programs.

Although this guide covers both the POSIX and Solaris threads implementations, most topics assume a POSIX threads interest. Information applying to only Solaris threads is covered in a special chapter.

To understand this guide, a reader must be familiar with

- A UNIX® SVR4 system—preferably the Solaris Operating Environment.
- The C programming language—multithreading is implemented through the `libthread` library
- The principles of concurrent programming (as opposed to sequential programming)—multithreading requires a different way of thinking about function interactions. Some books you might want to read are:
 - *Algorithms for Mutual Exclusion* by Michel Raynal (MIT Press, 1986)
 - *Concurrent Programming* by Alan Burns & Geoff Davies (Addison-Wesley, 1993)
 - *Distributed Algorithms and Protocols* by Michel Raynal (Wiley, 1988)
 - *Operating System Concepts* by Silberschatz, Peterson, & Galvin (Addison-Wesley, 1991)
 - *Principles of Concurrent Programming* by M. Ben-Ari (Prentice-Hall, 1982)

How This Guide Is Organized

Chapter 1 gives a structural overview of threads implementation in this release.

Chapter 2 discusses the general POSIX threads library routines, emphasizing creating a thread with default attributes.

Chapter 3 covers creating a thread with nondefault attributes.

Chapter 4 covers the threads library synchronization routines.

Chapter 5 discusses changes to the operating environment to support multithreading.

Chapter 6 covers multithreading safety issues.

Chapter 7 covers the basics of compiling and debugging multithreaded applications.

Chapter 8 covers the Solaris threads (as opposed to POSIX threads) interfaces.

Chapter 9 discusses issues that affect programmers writing multithreaded applications.

Appendix A shows how code can be designed for POSIX threads.

Appendix B shows an example of building a barrier in Solaris threads.

Accessing Sun Documentation Online

The docs.sun.comSM Web site enables you to access Sun technical documentation online. You can browse the docs.sun.com archive or search for a specific book title or subject. The URL is <http://docs.sun.com>.

What Typographic Conventions Mean

The following table describes the typographic changes used in this book.

TABLE P-1 Typographic Conventions

Typeface or Symbol	Meaning	Example
AaBbCc123	The names of commands, files, and directories; on-screen computer output	Edit your <code>.login</code> file. Use <code>ls -a</code> to list all files. <code>machine_name%</code> you have mail.
AaBbCc123	What you type, contrasted with on-screen computer output	<code>machine_name%</code> su Password:
<i>AaBbCc123</i>	Command-line placeholder: replace with a real name or value	To delete a file, type rm <i>filename</i> .
<i>AaBbCc123</i>	Book titles, new words, or terms, or words to be emphasized.	Read Chapter 6 in <i>User's Guide</i> . These are called <i>class</i> options. You must be <i>root</i> to do this.

Shell Prompts in Command Examples

The following table shows the default system prompt and superuser prompt for the C shell, Bourne shell, and Korn shell.

TABLE P-2 Shell Prompts

Shell	Prompt
C shell prompt	<code>machine_name%</code>
C shell superuser prompt	<code>machine_name#</code>
Bourne shell and Korn shell prompt	<code>\$</code>
Bourne shell and Korn shell superuser prompt	<code>#</code>

Covering Multithreading Basics

The word *multithreading* can be translated as *multiple threads of control* or *multiple flows of control*. While a traditional UNIX process always has contained and still does contain a single thread of control, multithreading (MT) separates a process into many execution threads, each of which runs independently.

Multithreading your code can

- Improve application responsiveness
- Use multiprocessors more efficiently
- Improve program structure
- Use fewer system resources

This chapter explains some multithreading terms, benefits, and concepts. If you are ready to start using multithreading, skip to Chapter 2.

- “Defining Multithreading Terms” on page 15
- “Meeting Multithreading Standards” on page 17
- “Benefiting From Multithreading” on page 17
- “Understanding Basic Multithreading Concepts” on page 18

Defining Multithreading Terms

Table 1-1 introduces some of the terms used in this book.

TABLE 1-1 Multithreading Terms

Term	Definition
Process	The UNIX environment (such as file descriptors, user ID, and so on) created with the <code>fork(2)</code> system call, which is set up to run a program.
Thread	A sequence of instructions executed within the context of a process.
pthread (POSIX threads)	A POSIX 1003.1c compliant threads interface.
Solaris threads	A Sun Microsystems™ threads interface that is not POSIX compliant. A predecessor of pthread.
Single-threaded	Restricting access to a single thread.
Multithreaded	Allowing access to two or more threads.
User- or Application-level threads	Threads managed by the threads library routines in user (as opposed to kernel) space.
Lightweight processes	Threads in the kernel that execute kernel code and system calls (also called LWPs).
Bound thread	A user-level thread that is permanently bound to one LWP.
Unbound thread	A user-level thread that is not necessarily bound to one LWP.
Attribute object	Contains opaque data types and related manipulation functions used to standardize some of the configurable aspects of POSIX threads, mutual exclusion locks (mutexes), and condition variables.
Mutual exclusion locks	Functions that lock and unlock access to shared data.
Condition variables	Functions that block threads until a change of state.
Read-write locks	Functions that allow multiple read-only access to shared data, but exclusive access for modification of that data.
Counting semaphore	A memory-based synchronization mechanism.
Parallelism	A condition that arises when at least two threads are <i>executing</i> simultaneously.
Concurrency	A condition that exists when at least two threads are <i>making progress</i> . A more generalized form of parallelism that can include time-slicing as a form of virtual parallelism.

Meeting Multithreading Standards

The concept of multithreaded programming goes back to at least the 1960s. Its development on UNIX systems began in the mid-1980s. While there is agreement about what multithreading is and the features necessary to support it, the interfaces used to implement multithreading have varied greatly.

For several years a group called POSIX (Portable Operating System Interface) 1003.4a has been working on standards for multithreaded programming. The standard has now been ratified. This Multithreaded Programming Guide is based on the POSIX standards: P1003.1b final draft 14 (realtime), and P1003.1c final draft 10 (multithreading).

This guide covers both POSIX threads (also called *pthread*s) and Solaris threads. Solaris threads were available in the Solaris 2.4 release, and are not functionally different from POSIX threads. However, because POSIX threads are more portable than Solaris threads, this guide covers multithreading from the POSIX perspective. Subjects specific to Solaris threads only are covered in the Chapter 8.

Benefiting From Multithreading

Improving Application Responsiveness

Any program in which many activities are not dependent upon each other can be redesigned so that each activity is defined as a thread. For example, the user of a multithreaded GUI does not have to wait for one activity to complete before starting another.

Using Multiprocessors Efficiently

Typically, applications that express concurrency requirements with threads need not take into account the number of available processors. The performance of the application improves transparently with additional processors.

Numerical algorithms and applications with a high degree of parallelism, such as matrix multiplications, can run much faster when implemented with threads on a multiprocessor.

Improving Program Structure

Many programs are more efficiently structured as multiple independent or semi-independent units of execution instead of as a single, monolithic thread. Multithreaded programs can be more adaptive to variations in user demands than single-threaded programs.

Using Fewer System Resources

Programs that use two or more processes that access common data through shared memory are applying more than one thread of control.

However, each process has a full address space and operating environment state. The cost of creating and maintaining this large amount of state information makes each process much more expensive than a thread in both time and space.

In addition, the inherent separation between processes can require a major effort by the programmer to communicate between the threads in different processes, or to synchronize their actions.

Combining Threads and RPC

By combining threads and a remote procedure call (RPC) package, you can exploit nonshared-memory multiprocessors (such as a collection of workstations). This combination distributes your application relatively easily and treats the collection of workstations as a multiprocessor.

For example, one thread might create child threads. Each of these children could then place a remote procedure call, invoking a procedure on another workstation. Although the original thread has merely created threads that are now running in parallel, this parallelism involves other computers.

Understanding Basic Multithreading Concepts

Concurrency and Parallelism

In a multithreaded process on a single processor, the processor can switch execution resources between threads, resulting in concurrent execution.

In the same multithreaded process in a shared-memory multiprocessor environment, each thread in the process can run on a separate processor at the same time, resulting in parallel execution.

When the process has fewer or as many threads as there are processors, the threads support system in conjunction with the operating environment ensure that each thread runs on a different processor.

For example, in a matrix multiplication that has the same number of threads and processors, each thread (and each processor) computes a row of the result.

Looking at Multithreading Structure

Traditional UNIX already supports the concept of threads—each process contains a single thread, so programming with multiple processes is programming with multiple threads. But a process is also an address space, and creating a process involves creating a new address space.

Creating a thread is less expensive when compared to creating a new process, because the newly created thread uses the current process address space. The time it takes to switch between threads is less than the time it takes to switch between processes, partly because switching between threads does not involve switching between address spaces.

Communicating between the threads of one process is simple because the threads share everything—address space, in particular. So, data produced by one thread is immediately available to all the other threads.

The interface to multithreading support is through a subroutine library, `libpthread` for POSIX threads, and `libthread` for Solaris threads. Multithreading provides flexibility by decoupling kernel-level and user-level resources.

User-Level Threads

Threads are the primary programming interface in multithreaded programming.¹ Threads are visible only from within the process, where they share all process resources like address space, open files, and so on. The following state is unique to each thread.

- Thread ID
- Register state (including PC and stack pointer)
- Stack
- Signal mask
- Priority

¹ User-level threads are so named to distinguish them from kernel-level threads, which are the concern of systems programmers only. Because this book is for application programmers, kernel-level threads are not discussed.

- Thread-private storage

Because threads share the process instructions and most of the process data, a change in shared data by one thread can be seen by the other threads in the process. When a thread needs to interact with other threads in the same process, it can do so without involving the operating environment.

By default, threads are lightweight. But, to get more control over a thread (for instance, to control scheduling policy more), the application can bind the thread. When an application binds threads to execution resources, the threads become kernel resources (see “System Scope (Bound Threads)” on page 22 for more information).

To summarize, user-level threads are:

- Inexpensive to create because they do not need to create their own address space.
- Fast to synchronize because synchronization is done at the application level, not at the kernel level.
- Managed by the threads library; either `libpthreads` or `libthread`.

Lightweight Processes

The threads library uses underlying threads of control called *lightweight processes* that are supported by the kernel. You can think of an LWP as a virtual CPU that executes code or system calls.

You usually do not need to concern yourself with LWPs to program with threads. The information here about LWPs is provided as background, so you can understand the differences in scheduling scope, described on “Process Scope (Unbound Threads)” on page 21.

Much as the `stdio` library routines such as `fopen()` and `fread()` use the `open()` and `read()` functions, the threads interface uses the LWP interface, and for many of the same reasons.

Lightweight processes (LWPs) bridge the user level and the kernel level. Each process contains one or more LWP, each of which runs one or more user threads. (See Figure 1-1.)

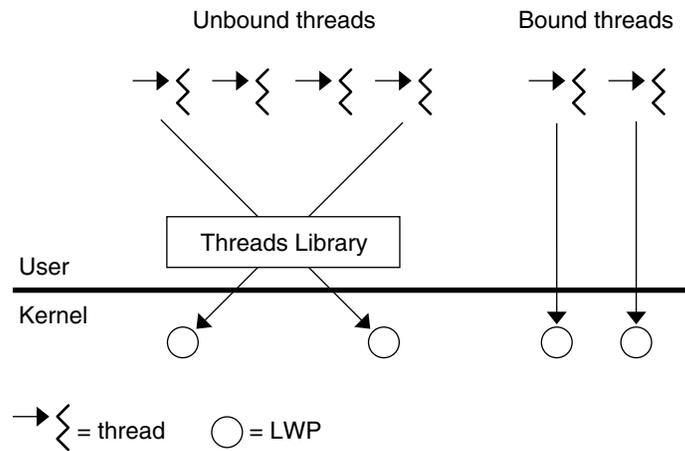


FIGURE 1-1 User-level Threads and Lightweight Processes

Each LWP is a kernel resource in a kernel pool, and is allocated and de-allocated to a thread on a per thread basis.

Scheduling

POSIX specifies three scheduling policies: first-in-first-out (`SCHED_FIFO`), round-robin (`SCHED_RR`), and custom (`SCHED_OTHER`). `SCHED_FIFO` is a queue-based scheduler with different queues for each priority level. `SCHED_RR` is like FIFO except that each thread has an execution time quota.

Both `SCHED_FIFO` and `SCHED_RR` are POSIX Realtime extensions. `SCHED_OTHER` is the default scheduling policy.

See “LWPs and Scheduling Classes” on page 142 for information about the `SCHED_OTHER` policy.

Two scheduling scopes are available: process scope for unbound threads and system scope for bound threads. Threads with differing scope states can coexist on the same system and even in the same process. In general, the scope sets the range in which the threads scheduling policy is in effect.

Process Scope (Unbound Threads)

`PTHREAD_SCOPE_PROCESS` threads are created as unbound threads. The association of these threads with LWPs is managed by the threads library.

In most cases, threads should be `PTHREAD_SCOPE_PROCESS`. These threads have no restriction to execute on a particular LWP, and are equivalent to Solaris thread created without the `THR_BOUND` flag. The threads library decides the association between individual threads and LWPs.

System Scope (Bound Threads)

`PTHREAD_SCOPE_SYSTEM` threads are created as bound threads. A bound thread is permanently attached to an LWP.

Each bound thread is bound to an LWP for the lifetime of the thread. This is equivalent to creating a Solaris thread in the `THR_BOUND` state. You can bind a thread to use special scheduling attributes with Realtime scheduling.

Note – In neither case, bound or unbound, can a thread be directly accessed by or moved to another process.

Cancellation

Thread cancellation allows a thread to terminate the execution of any other thread in the process. The target thread (the one being cancelled) can keep cancellation requests pending and can perform application-specific cleanup when it acts upon the cancellation notice.

The pthreads cancellation feature permits either asynchronous or deferred termination of a thread. Asynchronous cancellation can occur at any time; deferred cancellation can occur only at defined points. Deferred cancellation is the default type.

Synchronization

Synchronization allows you to control program flow and access to shared data for concurrently executing threads.

The four synchronization models are mutex locks, read/write locks, condition variables, and semaphores.

- Mutex locks allow only one thread at a time to execute a specific section of code, or to access specific data.
- Read/write locks permit concurrent reads and exclusive writes to a protected shared resource. To modify a resource, a thread must first acquire the exclusive write lock. An exclusive write lock is not permitted until all read locks have been released.

- Condition variables block threads until a particular condition is true.
- Counting semaphores typically coordinate access to resources. The count is the limit on how many threads can have access to a semaphore. When the count is reached, the semaphore blocks.

Using the 64-bit Architecture

For application developers, the major difference between the Solaris 64-bit and 32-bit operating environments is the C-language data type model used. The 64-bit data type uses the LP64 model where `longs` and pointers are 64-bits wide. All other fundamental data types remain the same as those of the 32-bit implementation. The 32-bit data type uses the ILP32 model where `ints`, `longs`, and pointers are 32-bits.

The following summary briefly describes the major features and considerations for using the 64-bit environment:

- Large Virtual Address Space

In the 64-bit environment, a process can have up to 64 bits of virtual address space, or 18 exabytes. This is 4 billion times the current 4 Gbyte maximum of a 32-bit process. Because of hardware restrictions, however, some platforms might not support the full 64 bits of address space.

Large address space increases the number of threads that can be created with the default stack size (1 megabyte on 32 bits, 2 megabytes on 64 bits). The number of threads with the default stack size is approximately 2000 threads on a 32-bit system and 8000 billion on a 64-bit system.

- Kernel Memory Readers

Because the kernel is an LP64 object that uses 64-bit data structures internally, existing 32-bit applications that use `libkvm`, `/dev/mem`, or `/dev/kmem` do not work properly and must be converted to 64-bit programs.

- `/proc` Restrictions

A 32-bit program that uses `/proc` is able to look at 32-bit processes, but is unable to understand a 64-bit process; the existing interfaces and data structures that describe the process are not large enough to contain the 64-bit quantities involved. Such programs must be recompiled as 64-bit programs to work for both 32-bit and 64-bit processes.

- 64-bit Libraries

32-bit applications are required to link with 32-bit libraries, and 64-bit applications are required to link with 64-bit libraries. With the exception of those libraries that have become obsolete, all of the system libraries are provided in both 32-bit and 64-bit versions. However, no 64-libraries are provided in static form.

- 64-bit Arithmetic

Though 64-bit arithmetic has long been available in previous 32-bit Solaris releases, the 64-bit implementation now provides full 64-bit machine registers for integer operations and parameter passing.

- Large Files

If an application requires only large file support, then it can remain 32-bit and use the Large Files interface. It is, however, recommended that the application be converted to 64-bit to take full advantage of 64-bit capabilities.

Basic Threads Programming

The Threads Library

This chapter introduces the basic threads programming routines from the POSIX threads library, `libpthread(3THR)`. This chapter covers *default threads*, or threads with default attribute values, which are the kind of threads that are most often used in multithreaded programming.

Chapter 3, explains how to create and use threads with nondefault attributes.

The POSIX (`libpthread`) routines introduced here have programming interfaces that are similar to the original (`libthread`) Solaris multithreading library.

The following brief roadmap directs you to the discussion of a particular task and its associated man page.

- “Create a Default Thread” on page 26
- “Wait for Thread Termination” on page 27
- “Detaching a Thread” on page 29
- “Create a Key for Thread-Specific Data” on page 30
- “Delete the Thread-Specific Data Key” on page 31
- “Set Thread-Specific Data” on page 32
- “Get Thread-Specific Data” on page 33
- “Get the Thread Identifier” on page 36
- “Compare Thread IDs” on page 36
- “Initializing Threads” on page 37
- “Yield Thread Execution” on page 37
- “Set the Thread Priority” on page 38
- “Get the Thread Priority” on page 39
- “Send a Signal to a Thread” on page 39
- “Access the Signal Mask of the Calling Thread” on page 40
- “Forking Safely” on page 41

- “Terminate a Thread” on page 41
- “Cancel a Thread” on page 44
- “Enable or Disable Cancellation” on page 45
- “Set Cancellation Type” on page 45
- “Create a Cancellation Point” on page 46
- “Push a Handler Onto the Stack” on page 47
- “Pull a Handler Off the Stack” on page 47

Create a Default Thread

When an attribute object is not specified, it is `NULL`, and the default thread is created with the following attributes:

- Unbound
- Nondetached
- With a default stack and stack size
- With a priority of zero

You can also create a default attribute object with `pthread_attr_init()`, and then use this attribute object to create a default thread. See the section “Initialize Attributes” on page 50 for details.

`pthread_create(3THR)`

Use `pthread_create(3THR)` to add a new thread of control to the current process.

```

Prototype:
int pthread_create(pthread_t *tid, const pthread_attr_t *tattr,
    void*(*start_routine)(void *), void *arg);

#include <pthread.h>

pthread_attr_t () tattr;
pthread_t tid;
extern void *start_routine(void *arg);
void *arg;
int ret;

/* default behavior*/
ret = pthread_create(&tid, NULL, start_routine, arg);

/* initialized with default attributes */
ret = pthread_attr_init(&tattr);
/* default behavior specified*/
ret = pthread_create(&tid, &tattr, start_routine, arg);

```

The `pthread_create()` function is called with `attr` having the necessary state behavior. `start_routine` is the function with which the new thread begins execution. When `start_routine` returns, the thread exits with the exit status set to the value returned by `start_routine` (see “`pthread_create(3THR)`” on page 26).

When `pthread_create()` is successful, the ID of the thread created is stored in the location referred to as *tid*.

Creating a thread using a NULL attribute argument has the same effect as using a default attribute; both create a default thread. When *tattr* is initialized, it acquires the default behavior.

Return Values

`pthread_create()` returns zero when it completes successfully. Any other return value indicates that an error occurred. When any of the following conditions are detected, `pthread_create()` fails and returns the corresponding value.

EAGAIN

A system limit is exceeded, such as when too many LWPs have been created.

EINVAL

The value of *tattr* is invalid.

Wait for Thread Termination

pthread_join(3THR)

Use `pthread_join(3THR)` to wait for a thread to terminate.

Prototype:

```
int pthread_join(thread_t tid, void **status);
```

```
#include <pthread.h>
```

```
pthread_t tid;
```

```
int ret;
```

```
void *status;
```

```
/* waiting to join thread "tid" with status */
```

```
ret = pthread_join(tid, &status);
```

```
/* waiting to join thread "tid" without status */
```

```
ret = pthread_join(tid, NULL);
```

The `pthread_join()` function blocks the calling thread until the specified thread terminates.

The specified thread must be in the current process and must not be detached. For information on thread detachment, see “Set Detach State” on page 52.

When *status* is not NULL, it points to a location that is set to the exit status of the terminated thread when `pthread_join()` returns successfully.

If multiple threads wait for the same thread to terminate, they all wait until the target thread terminates, then one thread returns successfully and the others fail with an error of `ESRCH`.

After `pthread_join()` returns, any data storage associated with the thread can be reclaimed by the application.

Return Values

`pthread_join()` returns zero when it completes successfully. Any other return value indicates that an error occurred. When any of the following conditions are detected, `pthread_join()` fails and returns the corresponding value.

`ESRCH`

tid is not a valid, undetached thread in the current process.

`EDEADLK`

A deadlock would exist, such as a thread waits for itself or thread A waits for thread B and thread B waits for thread A.

`EINVAL`

The value of *tid* is invalid.

Remember that `pthread_join()` works only for target threads that are nondetached. When there is no reason to synchronize with the termination of a particular thread, then that thread should be detached.

A Simple Threads Example

In Example 2-1, one thread executes the procedure at the top, creating a helper thread that executes the procedure `fetch()`, which involves a complicated database lookup and might take some time.

The main thread wants the results of the lookup but has other work to do in the meantime. So it does those other things and then waits for its helper to complete its job by executing `pthread_join()`.

An argument, *pbe*, to the new thread is passed as a stack parameter. This can be done here because the main thread waits for the spun-off thread to terminate. In general, though, it is better to use `malloc(3C)` to allocate storage from the heap instead of passing an address to thread stack storage, because this address might disappear or be reassigned if the thread terminated.

EXAMPLE 2-1 A Simple Threads Program

```
void mainline (...)  
{  
    struct phonebookentry *pbe;
```

EXAMPLE 2-1 A Simple Threads Program (Continued)

```
pthread_attr_t tattr;
pthread_t helper;
void *status;

pthread_create(&helper, NULL, fetch, &pbe);

    /* do something else for a while */

pthread_join(helper, &status);
/* it's now safe to use result */
}

void *fetch(struct phonebookentry *arg)
{
    struct phonebookentry *npbe;
    /* fetch value from a database */

    npbe = search (prog_name)
        if (npbe != NULL)
            *arg = *npbe;
    pthread_exit(0);
}

struct phonebookentry {
    char name[64];
    char phonenumber[32];
    char flags[16];
}
```

Detaching a Thread

pthread_detach(3THR)

pthread_detach(3THR) is an alternative to pthread_join(3THR) to reclaim storage for a thread that is created with a *detachstate* attribute set to PTHREAD_CREATE_JOINABLE.

Prototype:

```
int pthread_detach(pthread_t tid);
```

```
#include <pthread.h>
```

```
pthread_t tid;
int ret;
```

```
/* detach thread tid */
ret = pthread_detach(tid);
```

The `pthread_detach()` function is used to indicate to the implementation that storage for the thread *tid* can be reclaimed when the thread terminates. If *tid* has not terminated, `pthread_detach()` does not cause it to terminate. The effect of multiple `pthread_detach()` calls on the same target thread is unspecified.

Return Values

`pthread_detach()` returns zero when it completes successfully. Any other return value indicates that an error occurred. When any of the following conditions is detected, `pthread_detach()` fails and returns the corresponding value.

EINVAL

tid is not a valid thread.

ESRCH

tid is not a valid, undetached thread in the current process.

Create a Key for Thread-Specific Data

Single-threaded C programs have two basic classes of data—local data and global data. For multithreaded C programs a third class is added—*thread-specific data* (TSD). This is very much like global data, except that it is private to a thread.

Thread-specific data is maintained on a per-thread basis. TSD is the only way to define and refer to data that is private to a thread. Each thread-specific data item is associated with a *key* that is global to all threads in the process. Using the *key*, a thread can access a pointer (*void **) that is maintained per-thread.

pthread_key_create(3THR)

Use `pthread_key_create(3THR)` to allocate a *key* that is used to identify thread-specific data in a process. The *key* is global to all threads in the process, and all threads initially have the value NULL associated with the key when it is created.

Call `pthread_key_create()` once for each *key* before using the *key*. There is no implicit synchronization.

Once a *key* has been created, each thread can bind a value to the *key*. The values are specific to the threads and are maintained for each thread independently. The per-thread binding is deallocated when a thread terminates if the *key* was created with a destructor function.

Prototype:

```
int pthread_key_create(pthread_key_t *key,
                      void (*destructor) (void *));
```

```

#include <pthread.h>

pthread_key_t key;
int ret;

/* key create without destructor */
ret = pthread_key_create(&key, NULL);

/* key create with destructor */
ret = pthread_key_create(&key, destructor);

```

When `pthread_key_create()` returns successfully, the allocated key is stored in the location pointed to by *key*. The caller must ensure that the storage and access to this key are properly synchronized.

An optional destructor function, *destructor*, can be used to free stale storage. When a key has a non-NULL destructor function and the thread has a non-NULL value associated with that key, the destructor function is called with the current associated value when the thread exits. The order in which the destructor functions are called is unspecified.

Return Values

`pthread_key_create()` returns zero after completing successfully. Any other return value indicates that an error occurred. When any of the following conditions occur, `pthread_key_create()` fails and returns the corresponding value.

EAGAIN

The *key* name space is exhausted.

ENOMEM

Not enough virtual memory is available in this process to create a new key.

Delete the Thread-Specific Data Key

`pthread_key_delete(3THR)`

Use `pthread_key_delete(3THR)` to destroy an existing thread-specific data key. Any memory associated with the key can be freed because the key has been invalidated and will return an error if ever referenced. There is no comparable function in Solaris threads.

Prototype:

```

int pthread_key_delete(pthread_key_t key);

#include <pthread.h>

pthread_key_t key;

```

```
int ret;

/* key previously created */
ret = pthread_key_delete(key);
```

Once a *key* has been deleted, any reference to it with the `pthread_setspecific()` or `pthread_getspecific()` call yields undefined results.

It is the responsibility of the programmer to free any thread-specific resources before calling the delete function. This function does not invoke any of the destructors.

Return Values

`pthread_key_delete()` returns zero after completing successfully. Any other return value indicates that an error occurred. When the following condition occurs, `pthread_key_create()` fails and returns the corresponding value.

EINVAL
The *key* value is invalid.

Set Thread-Specific Data

`pthread_setspecific(3THR)`

Use `pthread_setspecific(3THR)` to set the thread-specific binding to the specified thread-specific data key.

```
Prototype:
int pthread_setspecific(pthread_key_t key, const void *value);

#include <pthread.h>

pthread_key_t key;
void *value;
int ret;

/* key previously created */
ret = pthread_setspecific(key, value);
```

Return Values

`pthread_setspecific()` returns zero after completing successfully. Any other return value indicates that an error occurred. When any of the following conditions occur, `pthread_setspecific()` fails and returns the corresponding value.

ENOMEM
Not enough virtual memory is available.

EINVAL
key is invalid.

Note – `pthread_setspecific()` does not free its storage. If a new binding is set, the existing binding must be freed; otherwise, a memory leak can occur.

Get Thread-Specific Data

`pthread_getspecific(3THR)`

Use `pthread_getspecific(3THR)` to get the calling thread's binding for *key*, and store it in the location pointed to by *value*.

```
Prototype:
void      *pthread_getspecific(pthread_key_t key);

#include <pthread.h>

pthread_key_t key;
void *value;

/* key previously created */
value = pthread_getspecific(key);
```

Return Values

No errors are returned.

Global and Private Thread-Specific Data Example

Example 2–2 shows an excerpt from a multithreaded program. This code is executed by any number of threads, but it has references to two global variables, *errno* and *mywindow*, that really should be references to items private to each thread.

EXAMPLE 2-2 Thread-Specific Data—Global but Private

```
body() {
    ...

    while (write(fd, buffer, size) == -1) {
        if (errno != EINTR) {
            fprintf(mywindow, "%s\n", strerror(errno));
            exit(1);
        }
    }
}
```

EXAMPLE 2-2 Thread-Specific Data—Global but Private (Continued)

```
    }  
    ...  
}
```

References to `errno` should get the system error code from the routine called by this thread, not by some other thread. So, references to `errno` by one thread refer to a different storage location than references to `errno` by other threads.

The `mywindow` variable is intended to refer to a `stdio` stream connected to a window that is private to the referring thread. So, as with `errno`, references to `mywindow` by one thread should refer to a different storage location (and, ultimately, a different window) than references to `mywindow` by other threads. The only difference here is that the threads library takes care of `errno`, but the programmer must somehow make this work for `mywindow`.

The next example shows how the references to `mywindow` work. The preprocessor converts references to `mywindow` into invocations of the `_mywindow()` procedure.

This routine in turn invokes `pthread_getspecific()`, passing it the `mywindow_key` global variable (it really is a global variable) and an output parameter, `win`, that receives the identity of this thread's window.

EXAMPLE 2-3 Turning Global References Into Private References

```
thread_key_t mywin_key;  
  
FILE *_mywindow(void) {  
    FILE *win;  
  
    win = pthread_getspecific(mywin_key);  
    return(win);  
}  
  
#define mywindow _mywindow()  
  
void routine_uses_win( FILE *win) {  
    ...  
}  
  
void thread_start(...) {  
    ...  
    make_mywin();  
    ...  
    routine_uses_win( mywindow )  
    ...  
}
```

The *mywin_key* variable identifies a class of variables for which each thread has its own private copy; that is, these variables are thread-specific data. Each thread calls `make_mywin()` to initialize its window and to arrange for its instance of *mywindow* to refer to it.

Once this routine is called, the thread can safely refer to *mywindow* and, after `_mywindow()`, the thread gets the reference to its private window. So, references to *mywindow* behave as if they were direct references to data private to the thread.

Example 2-4 shows how to set this up.

EXAMPLE 2-4 Initializing the Thread-Specific Data

```
void make_mywindow(void) {
    FILE **win;
    static pthread_once_t mykeycreated = PTHREAD_ONCE_INIT;

    pthread_once(&mykeycreated, mykeycreate);

    win = malloc(sizeof(*win));
    create_window(win, ...);

    pthread_setspecific(mywindow_key, win);
}

void mykeycreate(void) {
    pthread_key_create(&mywindow_key, free_key);
}

void free_key(void *win) {
    free(win);
}
```

First, get a unique value for the key, *mywin_key*. This key is used to identify the thread-specific class of data. So, the first thread to call `make_mywin()` eventually calls `pthread_key_create()`, which assigns to its first argument a unique *key*. The second argument is a destructor function that is used to deallocate a thread's instance of this thread-specific data item once the thread terminates.

The next step is to allocate the storage for the caller's instance of this thread-specific data item. Having allocated the storage, a call is made to the `create_window()` routine, which sets up a window for the thread and sets the storage pointed to by *win* to refer to it. Finally, a call is made to `pthread_setspecific()`, which associates the value contained in *win* (that is, the location of the storage containing the reference to the window) with the key.

After this, whenever this thread calls `pthread_getspecific()`, passing the global *key*, it gets the value that was associated with this key by this thread when it called `pthread_setspecific()`.

When a thread terminates, calls are made to the destructor functions that were set up in `pthread_key_create()`. Each destructor function is called only if the terminating thread established a value for the *key* by calling `pthread_setspecific()`.

Get the Thread Identifier

`pthread_self(3THR)`

Use `pthread_self(3THR)` to get the thread identifier of the calling thread.

```
Prototype:
pthread_t  pthread_self(void);

#include <pthread.h>

pthread_t  tid;

tid = pthread_self();
```

Return Values

`pthread_self()` returns the thread identifier of the calling thread.

Compare Thread IDs

`pthread_equal(3THR)`

Use `pthread_equal(3THR)` to compare the thread identification numbers of two threads.

```
Prototype:
int  pthread_equal(pthread_t tid1, pthread_t tid2);

#include <pthread.h>

pthread_t  tid1, tid2;
int  ret;

ret = pthread_equal(tid1, tid2);
```

Return Values

`pthread_equal()` returns a nonzero value when *tid1* and *tid2* are equal; otherwise, zero is returned. When either *tid1* or *tid2* is an invalid thread identification number, the result is unpredictable.

Initializing Threads

pthread_once(3THR)

Use `pthread_once(3THR)` to call an initialization routine the first time `pthread_once(3THR)` is called. Subsequent calls to `pthread_once()` have no effect.

```
Prototype:
int      pthread_once(pthread_once_t *once_control,
                      void (*init_routine)(void));

#include <pthread.h>

pthread_once_t once_control = PTHREAD_ONCE_INIT;
int ret;

ret = pthread_once(&once_control, init_routine);
```

The `once_control` parameter determines whether the associated initialization routine has been called.

Return Values

`pthread_once()` returns zero after completing successfully. Any other return value indicates that an error occurred. When the following condition occurs, `pthread_once()` fails and returns the corresponding value.

```
EINVAL
    once_control or init_routine is NULL.
```

Yield Thread Execution

sched_yield(3RT)

Use `sched_yield(3RT)` to cause the current thread to yield its execution in favor of another thread with the same or greater priority.

```
Prototype:
int      sched_yield(void);

#include <sched.h>

int ret;
```

```
ret = sched_yield();
```

Return Values

`sched_yield()` returns zero after completing successfully. Otherwise -1 is returned and *errno* is set to indicate the error condition.

ENOSYS

`sched_yield(3RT)` is not supported in this implementation.

Set the Thread Priority

`pthread_setschedparam(3THR)`

Use `pthread_setschedparam(3THR)` to modify the priority of an existing thread. This function has no effect on scheduling policy.

Prototype:

```
int pthread_setschedparam(pthread_t tid, int policy,
    const struct sched_param *param);
```

```
#include <pthread.h>
```

```
pthread_t tid;
int ret;
struct sched_param param;
int priority;
```

```
/* sched_priority will be the priority of the thread */
sched_param.sched_priority = priority;
```

```
/* only supported policy, others will result in ENOTSUP */
policy = SCHED_OTHER;
```

```
/* scheduling parameters of target thread */
ret = pthread_setschedparam(tid, policy, &param);
```

Return Values

`pthread_setschedparam()` returns zero after completing successfully. Any other return value indicates that an error occurred. When either of the following conditions occurs, the `pthread_setschedparam()` function fails and returns the corresponding value.

EINVAL

The value of the attribute being set is not valid.

ENOTSUP

An attempt was made to set the attribute to an unsupported value.

Get the Thread Priority

pthread_getschedparam(3THR)

pthread_getschedparam(3THR) gets the priority of the existing thread.

Prototype:

```
int pthread_getschedparam(pthread_t tid, int policy,
    struct schedparam *param);
```

```
#include <pthread.h>
```

```
pthread_t tid;
sched_param param;
int priority;
int policy;
int ret;
```

```
/* scheduling parameters of target thread */
ret = pthread_getschedparam (tid, &policy, &param);
```

```
/* sched_priority contains the priority of the thread */
priority = param.sched_priority;
```

Return Values

pthread_getschedparam() returns zero after completing successfully. Any other return value indicates that an error occurred. When the following condition occurs, the function fails and returns the corresponding value.

ESRCH

The value specified by *tid* does not refer to an existing thread.

Send a Signal to a Thread

pthread_kill(3THR)

Use pthread_kill(3THR) to send a signal to a thread.

Prototype:

```
int pthread_kill(pthread_t tid, int sig);
```

```

#include <pthread.h>
#include <signal.h>
int sig;
pthread_t tid;
int ret;

ret = pthread_kill(tid, sig);

```

`pthread_kill()` sends the signal *sig* to the thread specified by *tid*. *tid* must be a thread within the same process as the calling thread. The *sig* argument must be from the list given in `signal(5)`.

When *sig* is zero, error checking is performed but no signal is actually sent. This can be used to check the validity of *tid*.

Return Values

`pthread_kill()` returns zero after completing successfully. Any other return value indicates that an error occurred. When either of the following conditions occurs, `pthread_kill()` fails and returns the corresponding value.

EINVAL
sig is not a valid signal number.

ESRCH
tid cannot be found in the current process.

Access the Signal Mask of the Calling Thread

pthread_sigmask(3THR)

Use `pthread_sigmask(3THR)` to change or examine the signal mask of the calling thread.

```

Prototype:
int pthread_sigmask(int how, const sigset_t *new, sigset_t *old);

#include <pthread.h>
#include <signal.h>
int ret;
sigset_t old, new;

ret = pthread_sigmask(SIG_SETMASK, &new, &old); /* set new mask */
ret = pthread_sigmask(SIG_BLOCK, &new, &old); /* blocking mask */
ret = pthread_sigmask(SIG_UNBLOCK, &new, &old); /* unblocking */

```

how determines how the signal set is changed. It can have one of the following values:

- `SIG_BLOCK`—Add *new* to the current signal mask, where *new* indicates the set of signals to block.
- `SIG_UNBLOCK`—Delete *new* from the current signal mask, where *new* indicates the set of signals to unblock.
- `SIG_SETMASK`—Replace the current signal mask with *new*, where *new* indicates the new signal mask.

When the value of *new* is `NULL`, the value of *how* is not significant and the signal mask of the thread is unchanged. So, to inquire about currently blocked signals, assign a `NULL` value to the *new* argument.

The *old* variable points to the space where the previous signal mask is stored, unless it is `NULL`.

Return Values

`pthread_sigmask()` returns zero when it completes successfully. Any other return value indicates that an error occurred. When the following condition occurs, `pthread_sigmask()` fails and returns the corresponding value.

`EINVAL`

The value of *how* is not defined.

Forking Safely

`pthread_atfork(3THR)`

See the discussion about `pthread_atfork(3THR)` in “The Solution—`pthread_atfork(3THR)`” on page 138.

Prototype:

```
int pthread_atfork(void (*prepare) (void), void (*parent) (void),
                  void (*child) (void) );
```

Terminate a Thread

`pthread_exit(3THR)`

Use `pthread_exit(3THR)` to terminate a thread.

```

Prototype:
void      pthread_exit(void *status);

#include <pthread.h>

void *status;

pthread_exit(status); /* exit with status */

```

The `pthread_exit()` function terminates the calling thread. All thread-specific data bindings are released. If the calling thread is not detached, then the thread's ID and the exit status specified by *status* are retained until the thread is waited for via `pthread_join()`. Otherwise, *status* is ignored and the thread's ID can be reclaimed immediately. For information on thread detachment, see "Set Detach State" on page 52.

Return Values

The calling thread terminates with its exit status set to the contents of *status*.

Finishing Up

A thread can terminate its execution in the following ways:

- By returning from its first (outermost) procedure, the thread's start routine; see `pthread_create(3THR)`
- By calling `pthread_exit()`, supplying an exit status
- By termination with POSIX cancel functions; see `pthread_cancel()`

The default behavior of a thread is to linger until some other thread has acknowledged its demise by "joining" with it. This is the same as the default `pthread_create()` attribute being `nondetached`; see `pthread_detach(3THR)`. The result of the join is that the joining thread picks up the exit status of the dying thread and the dying thread vanishes.

An important special case arises when the initial thread — the one calling `main()`, — returns from calling `main()` or calls `exit(3C)`. This action causes the entire process to be terminated, along with all its threads. So take care to ensure that the initial thread does not return from `main()` prematurely.

Note that when the main thread merely calls `pthread_exit(3THR)`, it terminates only itself—the other threads in the process, as well as the process, continue to exist. (The process terminates when all threads terminate.)

Cancellation

Cancellation allows a thread to terminate the execution of any other thread, or all threads, in the process. Cancellation is an option when all further operations of a related set of threads are undesirable or unnecessary.

One example of thread cancellation is an asynchronously generated cancel condition, such as, when a user requesting to close or exit some running application. Another example is the completion of a task undertaken by a number of threads. One of the threads might ultimately complete the task while the others continue to operate. Since they are serving no purpose at that point, they all should be cancelled.

There are dangers in performing cancellations. Most deal with properly restoring invariants and freeing shared resources. A thread that is cancelled without care might leave a mutex in a locked state, leading to a deadlock. Or it might leave a region of memory allocated with no way to identify it and therefore no way to free it.

The `pthread` library specifies a cancellation interface that permits or forbids cancellation programmatically. The library defines the set of points at which cancellation can occur (*cancellation points*). It also allows the scope of *cancellation handlers*, which provide clean up services, to be defined so that they are sure to operate when and where intended.

Placement of cancellation points and the effects of cancellation handlers must be based on an understanding of the application. A mutex is explicitly not a cancellation point and should be held only the minimal essential time.

Limit regions of asynchronous cancellation to sequences with no external dependencies that could result in dangling resources or unresolved state conditions. Take care to restore cancellation state when returning from some alternate, nested cancellation state. The interface provides features to facilitate restoration: `pthread_setcancelstate(3THR)` preserves the current cancel state in a referenced variable; `pthread_setcanceltype(3THR)` preserves the current cancel type in the same way.

Cancellations can occur under three different circumstances:

- Asynchronously
- At various points in the execution sequence as defined by the standard
- At discrete points specified by the application

By default, cancellation can occur only at well-defined points as defined by the POSIX standard.

In all cases, take care that resources and state are restored to a condition consistent with the point of origin.

Cancellation Points

Be careful to cancel a thread only when cancellation is safe. The `pthread` standard specifies several cancellation points, including:

- Programmatically establish a thread cancellation point through a `pthread_testcancel(3THR)` call.
- Threads waiting for the occurrence of a particular condition in `pthread_cond_wait(3THR)` or `pthread_cond_timedwait(3THR)`.
- Threads waiting for termination of another thread in `pthread_join(3THR)`.
- Threads blocked on `sigwait(2)`.
- Some standard library calls. In general, these are functions in which threads can block; see the man page `cancellation(3THR)` for a list.

Cancellation is enabled by default. At times you might want an application to disable cancellation. This has the result of deferring all cancellation requests until they are enabled again.

See “`pthread_setcancelstate(3THR)`” on page 45 for information about disabling cancellation.

Cancel a Thread

`pthread_cancel(3THR)`

Use `pthread_cancel(3THR)` to cancel a thread.

Prototype:

```
int    pthread_cancel(pthread_t thread);  
  
#include <pthread.h>  
  
pthread_t thread;  
int ret;  
  
ret = pthread_cancel(thread);
```

How the cancellation request is treated depends on the state of the target thread. Two functions, `pthread_setcancelstate(3THR)` and `pthread_setcanceltype(3THR)`, determine that state.

Return Values

`pthread_cancel()` returns zero after completing successfully. Any other return value indicates that an error occurred. When the following condition occurs, the function fails and returns the corresponding value.

ESRCH

No thread could be found corresponding to that specified by the given thread ID.

Enable or Disable Cancellation

pthread_setcancelstate(3THR)

Use `pthread_setcancelstate(3THR)` to enable or disable thread cancellation. When a thread is created, thread cancellation is enabled by default.

Prototype:

```
int    pthread_setcancelstate(int state, int *oldstate);
#include <pthread.h>

int    oldstate;
int    ret;

/* enabled */
ret = pthread_setcancelstate(PTHREAD_CANCEL_ENABLE, &oldstate);

/* disabled */
ret = pthread_setcancelstate(PTHREAD_CANCEL_DISABLE, &oldstate);
```

Return Values

`pthread_setcancelstate()` returns zero after completing successfully. Any other return value indicates that an error occurred. When the following condition occurs, the `pthread_setcancelstate()` function fails and returns the corresponding value.

EINVAL

The state is not `PTHREAD_CANCEL_ENABLE` or `PTHREAD_CANCEL_DISABLE`.

Set Cancellation Type

pthread_setcanceltype(3THR)

Use `pthread_setcanceltype(3THR)` to set the cancellation type to either deferred or asynchronous mode. When a thread is created, the cancellation type is set to deferred mode by default. In deferred mode, the thread can be cancelled only at cancellation points. In asynchronous mode, a thread can be cancelled at any point during its execution. Using asynchronous mode is discouraged.

Prototype:

```
int pthread_setcanceltype(int type, int *oldtype);  
#include <pthread.h>  
  
int oldtype;  
int ret;  
  
/* deferred mode */  
ret = pthread_setcanceltype(PTHREAD_CANCEL_DEFERRED, &oldtype);  
  
/* async mode*/  
ret = pthread_setcanceltype(PTHREAD_CANCEL_ASYNC, &oldtype);
```

Return Values

`pthread_setcanceltype()` returns zero after completing successfully. Any other return value indicates that an error occurred. When the following condition occurs, the function fails and returns the corresponding value.

EINVAL

The type is not `PTHREAD_CANCEL_DEFERRED` or `PTHREAD_CANCEL_ASYNC`.

Create a Cancellation Point

`pthread_testcancel(3THR)`

Use `pthread_testcancel(3THR)` to establish a cancellation point for a thread.

Prototype:

```
void pthread_testcancel(void);  
#include <pthread.h>  
  
pthread_testcancel();
```

The `pthread_testcancel()` function is effective when thread cancellation is enabled and in deferred mode. Calling this function while cancellation is disabled has no effect.

Be careful to insert `pthread_testcancel()` only in sequences where it is safe to cancel a thread. In addition to programmatically establishing cancellation points through the `pthread_testcancel()` call, the pthreads standard specifies several cancellation points. See "Cancellation Points" on page 44 for more details.

There is no return value.

Push a Handler Onto the Stack

Use cleanup handlers to restore conditions to a state consistent with that at the point of origin, such as cleaning up allocated resources and restoring invariants. Use the `pthread_cleanup_push(3THR)` and `pthread_cleanup_pop(3THR)` functions to manage the handlers.

Cleanup handlers are pushed and popped in the same lexical scope of a program. They should always match; otherwise compiler errors will be generated.

`pthread_cleanup_push(3THR)`

Use `pthread_cleanup_push(3THR)` to push a cleanup handler onto a cleanup stack (LIFO).

Prototype:

```
void pthread_cleanup_push(void(*routine)(void*), void *args);  
  
#include <pthread.h>  
  
/* push the handler "routine" on cleanup stack */  
pthread_cleanup_push (routine, arg);
```

Pull a Handler Off the Stack

`pthread_cleanup_pop(3THR)`

Use `pthread_cleanup_pop(3THR)` to pull the cleanup handler off the cleanup stack.

A nonzero argument in the pop function removes the handler from the stack and executes it. An argument of zero pops the handler without executing it.

`pthread_cleanup_pop()` is effectively called with a nonzero argument if a thread either explicitly or implicitly calls `pthread_exit(3THR)` or if the thread accepts a cancel request.

Prototype:

```
void pthread_cleanup_pop(int execute);  
  
#include <pthread.h>  
  
/* pop the "func" out of cleanup stack and execute "func" */
```

```
pthread_cleanup_pop (1);  
  
/* pop the "func" and DONT execute "func" */  
pthread_cleanup_pop (0);
```

There are no return values.

Thread Create Attributes

The previous chapter covered the basics of threads creation using default attributes. This chapter discusses setting attributes at thread creation time.

Note that only pthreads uses attributes and cancellation, so the API covered in this chapter is for POSIX threads only. Otherwise, the *functionality* for Solaris threads and pthreads is largely the same. (See Chapter 9, Programming With Solaris Threads, for more information about similarities and differences.)

- “Initialize Attributes” on page 50
- “Destroy Attributes” on page 52
- “Set Detach State” on page 52
- “Get Detach State” on page 53
- “Set Stack Guard Size” on page 54
- “Get Stack Guard Size” on page 55
- “Set Scope” on page 56
- “Get Scope” on page 57
- “Set Thread Concurrency Level” on page 57
- “Get Thread Concurrency Level” on page 58
- “Set Scheduling Policy” on page 58
- “Get Scheduling Policy” on page 59
- “Set Inherited Scheduling Policy” on page 60
- “Get Inherited Scheduling Policy” on page 61
- “Set Scheduling Parameters” on page 61
- “Get Scheduling Parameters” on page 62
- “Set Stack Size” on page 64
- “Get Stack Size” on page 65
- “Set Stack Address” on page 67
- “Get Stack Address” on page 68

Attributes

Attributes are a way to specify behavior that is different from the default. When a thread is created with `pthread_create(3THR)` or when a synchronization variable is initialized, an attribute object can be specified. The defaults are usually sufficient.

An attribute object is opaque, and cannot be directly modified by assignments. A set of functions is provided to initialize, configure, and destroy each object type.

Once an attribute is initialized and configured, it has process-wide scope. The suggested method for using attributes is to configure all required state specifications at one time in the early stages of program execution. The appropriate attribute object can then be referred to as needed.

Using attribute objects has two primary advantages.

- First, it adds to code portability.

Even though supported attributes might vary between implementations, you need not modify function calls that create thread entities because the attribute object is hidden from the interface.

If the target port supports attributes that are not found in the current port, provision must be made to manage the new attributes. This is an easy porting task though, because attribute objects need only be initialized once in a well-defined location.

- Second, state specification in an application is simplified.

As an example, consider that several sets of threads might exist within a process, each providing a separate service, and each with its own state requirements.

At some point in the early stages of the application, a thread attribute object can be initialized for each set. All future thread creations will then refer to the attribute object initialized for that type of thread. The initialization phase is simple and localized, and any future modifications can be made quickly and reliably.

Attribute objects require attention at process exit time. When the object is initialized, memory is allocated for it. This memory must be returned to the system. The `pthread` standard provides function calls to destroy attribute objects.

Initialize Attributes

`pthread_attr_init(3THR)`

Use `pthread_attr_init(3THR)` to initialize object attributes to their default values. The storage is allocated by the thread system during execution.

Prototype:

```
int pthread_attr_init(pthread_attr_t *tattr);  
  
#include <pthread.h>  
  
pthread_attr_t tattr;  
int ret;  
  
/* initialize an attribute to the default value */  
ret = pthread_attr_init(&tattr);
```

Table 3–1 shows the default values for attributes (*tattr*).

TABLE 3–1 Default Attribute Values for *tattr*

Attribute	Value	Result
<i>scope</i>	PTHREAD_SCOPE_PROCESS	New thread is unbound – not permanently attached to LWP.
<i>detachstate</i>	PTHREAD_CREATE_JOINABLE	Exit status and thread are preserved after the thread terminates.
<i>stackaddr</i>	NULL	New thread has system-allocated stack address.
<i>stacksize</i>	0	New thread has system-defined stack size.
<i>priority</i>	0	New thread has priority 0.
<i>inheritsched</i>	PTHREAD_EXPLICIT_SCHED	New thread does not inherit parent thread scheduling priority.
<i>schedpolicy</i>	SCHED_OTHER	New thread uses Solaris-defined fixed priorities for synchronization object contention; threads run until preempted or until they block or yield.

Return Values

`pthread_attr_init()` returns zero after completing successfully. Any other return value indicates that an error occurred. If the following condition occurs, the function fails and returns the corresponding value.

ENOMEM

Returned when there is not enough memory to initialize the thread attributes object.

Destroy Attributes

`pthread_attr_destroy(3THR)`

Use `pthread_attr_destroy(3THR)` to remove the storage allocated during initialization. The attribute object becomes invalid.

Prototype:

```
int pthread_attr_destroy(pthread_attr_t *tattr);
```

```
#include <pthread.h>
```

```
pthread_attr_t tattr;  
int ret;
```

```
/* destroy an attribute */  
ret = pthread_attr_destroy(&tattr);
```

Return Values

`pthread_attr_destroy()` returns zero after completing successfully. Any other return value indicates that an error occurred. If the following condition occurs, the function fails and returns the corresponding value.

EINVAL

Indicates that the value of *tattr* was not valid.

Set Detach State

`pthread_attr_setdetachstate(3THR)`

When a thread is created detached (`PTHREAD_CREATE_DETACHED`), its thread *ID* and other resources can be reused as soon as the thread terminates. Use `pthread_attr_setdetachstate(3THR)` when the calling thread does not want to wait for the thread to terminate.

When a thread is created nondetached (`PTHREAD_CREATE_JOINABLE`), it is assumed that you will be waiting for it. That is, it is assumed that you will be executing a `pthread_join()` on the thread.

Whether a thread is created detached or nondetached, the process does not exit until all threads have exited. See “Finishing Up” on page 42 for a discussion of process termination caused by premature exit from `main()`.

Prototype:

```
int pthread_attr_setdetachstate(pthread_attr_t *tattr, int detachstate);
```

```

#include <pthread.h>

pthread_attr_t tattr;
int ret;

/* set the thread detach state */
ret = pthread_attr_setdetachstate(&tattr, PTHREAD_CREATE_DETACHED);

```

Note – When there is no explicit synchronization to prevent it, a newly created, detached thread can die and have its thread ID reassigned to another new thread before its creator returns from `pthread_create()`.

For nondetached (`PTHREAD_CREATE_JOINABLE`) threads, it is very important that some thread join with it after it terminates—otherwise the resources of that thread are not released for use by new threads. This commonly results in a memory leak. So when you do not want a thread to be joined, create it as a detached thread.

EXAMPLE 3-1 Creating a Detached Thread

```

#include <pthread.h>

pthread_attr_t tattr;
pthread_t tid;
void *start_routine;
void arg;
int ret;

/* initialized with default attributes */
ret = pthread_attr_init(&tattr);
ret = pthread_attr_setdetachstate(&tattr, PTHREAD_CREATE_DETACHED);
ret = pthread_create(&tid, &tattr, start_routine, arg);

```

Return Values

`pthread_attr_setdetachstate()` returns zero after completing successfully. Any other return value indicates that an error occurred. If the following condition occurs, the function fails and returns the corresponding value.

`EINVAL`

Indicates that the value of *detachstate* or *tattr* was not valid.

Get Detach State

`pthread_attr_getdetachstate(3THR)`

Use `pthread_attr_getdetachstate(3THR)` to retrieve the thread create state, which can be either detached or joined.

Prototype:

```
int pthread_attr_getdetachstate(const pthread_attr_t *tattr,
                               int *detachstate);

#include <pthread.h>

pthread_attr_t tattr;
int detachstate;
int ret;

/* get detachstate of thread */
ret = pthread_attr_getdetachstate (&tattr, &detachstate);
```

Return Values

`pthread_attr_getdetachstate()` returns zero after completing successfully. Any other return value indicates that an error occurred. If the following condition occurs, the function fails and returns the corresponding value.

EINVAL

Indicates that the value of *detachstate* is NULL or *tattr* is invalid.

Set Stack Guard Size

pthread_attr_setguardsize(3THR)

`pthread_attr_setguardsize(3THR)` sets the *guardsize* of the *attr* object.

The *guardsize* argument provides protection against overflow of the stack pointer. If a thread's stack is created with guard protection, the implementation allocates extra memory at the overflow end of the stack as a buffer against stack overflow of the stack pointer. If an application overflows into this buffer an error results (possibly in a SIGSEGV signal being delivered to the thread).

The *guardsize* attribute is provided to the application for two reasons:

1. Overflow protection can potentially result in wasted system resources. An application that creates a large number of threads, and knows its threads will never overflow their stack, can save system resources by turning off guard areas.
2. When threads allocate large data structures on stack, a large guard area may be needed to detect stack overflow.

If *guardsize* is zero, a guard area will not be provided for threads created with *attr*. If *guardsize* is greater than zero, a guard area of at least size *guardsize* bytes is provided for each thread created with *attr*. By default, a thread has an implementation-defined, non-zero guard area.

A conforming implementation is permitted to round up the value contained in *guardsize* to a multiple of the configurable system variable `PAGESIZE` (see `PAGESIZE` in `sys/mman.h`). If an implementation rounds up the value of *guardsize* to a multiple of `PAGESIZE`, a call to `pthread_attr_getguardsize()` specifying *attr* will store, in *guardsize*, the guard size specified in the previous call to `pthread_attr_setguardsize()`.

```
#include <pthread.h>
```

```
int pthread_attr_setguardsize(pthread_attr_t *attr, size_t guardsize);
```

Return Value

`pthread_attr_setguardsize()` fails if:

`EINVAL`

The argument *attr* is invalid, the argument *guardsize* is invalid, or the argument *guardsize* contains an invalid value.

Get Stack Guard Size

`pthread_attr_getguardsize(3THR)`

`pthread_attr_getguardsize(3THR)` gets the *guardsize* of the *attr* object.

A conforming implementation is permitted to round up the value contained in *guardsize* to a multiple of the configurable system variable `PAGESIZE` (see `PAGESIZE` in `sys/mman.h`). If an implementation rounds up the value of *guardsize* to a multiple of `PAGESIZE`, a call to `pthread_attr_getguardsize()` specifying *attr* will store, in *guardsize*, the guard size specified in the previous call to `pthread_attr_setguardsize()`.

```
#include <pthread.h>
```

```
int pthread_attr_getguardsize(const pthread_attr_t *attr,  
                             size_t *guardsize);
```

Return Value

`pthread_attr_getguardsize()` fails if:

`EINVAL`

The argument *attr* is invalid, the argument *guardsize* is invalid, or the argument *guardsize* contains an invalid value.

Set Scope

pthread_attr_setscope(3THR)

Use `pthread_attr_setscope(3THR)` to create a bound thread (`PTHREAD_SCOPE_SYSTEM`) or an unbound thread (`PTHREAD_SCOPE_PROCESS`).

Note – Both thread types are accessible only within a given process.

Prototype:

```
int    pthread_attr_setscope(pthread_attr_t *tattr, int scope);
#include <pthread.h>

pthread_attr_t tattr;
int ret;

/* bound thread */
ret = pthread_attr_setscope(&tattr, PTHREAD_SCOPE_SYSTEM);

/* unbound thread */
ret = pthread_attr_setscope(&tattr, PTHREAD_SCOPE_PROCESS);
```

Notice that there are three function calls in this example: one to initialize the attributes, one to set any variations from the default attributes, and one to create the pthreads.

```
#include <pthread.h>

pthread_attr_t attr;
pthread_t tid;
void start_routine;
void arg;
int ret;

/* initialized with default attributes */
ret = pthread_attr_init (&tattr);

/* BOUND behavior */
ret = pthread_attr_setscope(&tattr, PTHREAD_SCOPE_SYSTEM);
ret = pthread_create (&tid, &tattr, start_routine, arg);
```

Return Values

`pthread_attr_setscope()` returns zero after completing *successfully*. Any other return value indicates that an error occurred. If the following conditions occur, the function fails and returns the corresponding value.

EINVAL

An attempt was made to set *tattr* to a value that is not valid.

Get Scope

pthread_attr_getscope(3THR)

Use `pthread_attr_getscope(3THR)` to retrieve the thread scope, which indicates whether the thread is bound or unbound.

Prototype:

```
int    pthread_attr_getscope(pthread_attr_t *tattr, int *scope);

#include <pthread.h>

pthread_attr_t tattr;
int scope;
int ret;

/* get scope of thread */
ret = pthread_attr_getscope(&tattr, &scope);
```

Return Values

`pthread_attr_getscope()` returns zero after completing successfully. Any other return value indicates that an error occurred. If the following condition occurs, the function fails and returns the corresponding value.

EINVAL

The value of *scope* is NULL or *tattr* is invalid.

Set Thread Concurrency Level

pthread_setconcurrency(3THR)

`pthread_setconcurrency(3THR)` is provided for standards compliance. It is used by an application to inform the threads library of its desired concurrency level. For the threads implementation introduced in Solaris 9, this interface has no effect; all runnable threads are attached to LWPs.

```
#include <pthread.h>

int pthread_setconcurrency(int new_level);
```

Return Value

`pthread_setconcurrency()` fails if:

`EINVAL`

The value specified by *new_level* is negative.

`EAGAIN`

The value specified by *new_level* would cause a system resource to be exceeded.

Get Thread Concurrency Level

`pthread_getconcurrency(3THR)`

`pthread_getconcurrency(3THR)` returns the value set by a previous call to `pthread_setconcurrency()`. If the `pthread_setconcurrency()` function was not previously called, `pthread_getconcurrency()` returns zero.

```
#include <pthread.h>
```

```
int pthread_getconcurrency(void);
```

Return Value

`pthread_getconcurrency()` always returns the concurrency level set by a previous call to `pthread_setconcurrency()`. If `pthread_setconcurrency()` has never been called, `pthread_getconcurrency()` returns zero.

Set Scheduling Policy

`pthread_attr_setschedpolicy(3THR)`

Use `pthread_attr_setschedpolicy(3THR)` to set the scheduling policy. The POSIX standard specifies scheduling policy attributes of `SCHED_FIFO` (first-in-first-out), `SCHED_RR` (round-robin), or `SCHED_OTHER` (an implementation-defined method).

- `SCHED_FIFO`

First-In-First-Out; threads whose contention scope is system (`PTHREAD_SCOPE_SYSTEM`) are in real-time (RT) scheduling class if the calling process has an effective user id of 0. These threads, if not preempted by a higher priority thread, will proceed until they yield or block. `SCHED_FIFO` for threads

that have a contention scope of process (`PTHREAD_SCOPE_PROCESS`) or whose calling process does not have an effective user id of 0 is based on the TS scheduling class.

■ `SCHED_RR`

Round-Robin; threads whose contention scope is system (`PTHREAD_SCOPE_SYSTEM`) are in real-time (RT) scheduling class if the calling process has an effective user id of 0. These threads, if not preempted by a higher priority thread, and if they do not yield or block, will execute for a time period determined by the system. `SCHED_RR` for threads that have a contention scope of process (`PTHREAD_SCOPE_PROCESS`) or whose calling process does not have an effective user id of 0 is based on the TS scheduling class.

`SCHED_FIFO` and `SCHED_RR` are optional in POSIX, and are supported for real time bound threads only.

For a discussion of scheduling, see the section “Scheduling” on page 21.

Prototype:

```
int    pthread_attr_setschedpolicy(pthread_attr_t *tattr, int policy);
#include <pthread.h>

pthread_attr_t tattr;
int policy;
int ret;

/* set the scheduling policy to SCHED_OTHER */
ret = pthread_attr_setschedpolicy(&tattr, SCHED_OTHER);
```

Return Values

`pthread_attr_setschedpolicy()` returns zero after completing successfully. Any other return value indicates that an error occurred. When either of the following conditions occurs, the function fails and returns the corresponding value.

`EINVAL`

An attempt was made to set *tattr* to a value that is not valid.

`ENOTSUP`

An attempt was made to set the attribute to an unsupported value.

Get Scheduling Policy

`pthread_attr_getschedpolicy(3THR)`

Use `pthread_attr_getschedpolicy(3THR)` to retrieve the scheduling policy.

Prototype:

```
int pthread_attr_getschedpolicy(pthread_attr_t *tattr, int *policy);

#include <pthread.h>

pthread_attr_t tattr;
int policy;
int ret;

/* get scheduling policy of thread */
ret = pthread_attr_getschedpolicy (&tattr, &policy);
```

Return Values

`pthread_attr_getschedpolicy()` returns zero after completing successfully. Any other return value indicates that an error occurred. If the following condition occurs, the function fails and returns the corresponding value.

EINVAL

The parameter *policy* is NULL or *tattr* is invalid.

Set Inherited Scheduling Policy

`pthread_attr_setinheritsched(3THR)`

Use `pthread_attr_setinheritsched(3THR)` to set the inherited scheduling policy.

An *inherit* value of `PTHREAD_INHERIT_SCHED` means that the scheduling policies defined in the creating thread are to be used, and any scheduling attributes defined in the `pthread_create()` call are to be ignored. If `PTHREAD_EXPLICIT_SCHED` (the default) is used, the attributes from the `pthread_create()` call are to be used.

Prototype:

```
int pthread_attr_setinheritsched(pthread_attr_t *tattr, int inherit);

#include <pthread.h>

pthread_attr_t tattr;
int inherit;
int ret;

/* use the current scheduling policy */
ret = pthread_attr_setinheritsched(&tattr, PTHREAD_EXPLICIT_SCHED);
```

Return Values

`pthread_attr_setinheritsched()` returns zero after completing successfully. Any other return value indicates that an error occurred. When either of the following conditions occurs, the function fails and returns the corresponding value.

EINVAL

An attempt was made to set *tattr* to a value that is not valid.

ENOTSUP

An attempt was made to set the attribute to an unsupported value.

Get Inherited Scheduling Policy

`pthread_attr_getinheritsched(3THR)`

`pthread_attr_getinheritsched(3THR)` returns the scheduling policy set by `pthread_attr_setinheritsched()`.

Prototype:

```
int    pthread_attr_getinheritsched(pthread_attr_t *tattr, int *inherit);
#include <pthread.h>

pthread_attr_t tattr;
int inherit;
int ret;

/* get scheduling policy and priority of the creating thread */
ret = pthread_attr_getinheritsched (&tattr, &inherit);
```

Return Values

`pthread_attr_getinheritsched()` returns zero after completing successfully. Any other return value indicates that an error occurred. If the following condition occurs, the function fails and returns the corresponding value.

EINVAL

The parameter *inherit* is NULL or *tattr* is invalid.

Set Scheduling Parameters

`pthread_attr_setschedparam(3THR)`

`pthread_attr_setschedparam(3THR)` sets the scheduling parameters.

Scheduling parameters are defined in the `param` structure; only priority is supported. Newly created threads run with this priority.

Prototype:

```
int pthread_attr_setschedparam(pthread_attr_t *tattr,
                               const struct sched_param *param);

#include <pthread.h>

pthread_attr_t tattr;
int newprio;
sched_param param;
newprio = 30;

/* set the priority; others are unchanged */
param.sched_priority = newprio;

/* set the new scheduling param */
ret = pthread_attr_setschedparam (&tattr, &param);
```

Return Values

`pthread_attr_setschedparam()` returns zero after completing successfully. Any other return value indicates that an error occurred. If the following conditions occur, the function fails and returns the corresponding value.

EINVAL

The value of *param* is NULL or *tattr* is invalid.

You can manage pthreads priority two ways. You can set the priority attribute before creating a child thread, or you can change the priority of the parent thread and then change it back.

Get Scheduling Parameters

`pthread_attr_getschedparam(3THR)`

`pthread_attr_getschedparam(3THR)` returns the scheduling parameters defined by `pthread_attr_setschedparam()`.

Prototype:

```
int pthread_attr_getschedparam(pthread_attr_t *tattr,
                               const struct sched_param *param);

#include <pthread.h>

pthread_attr_t attr;
```

```

struct sched_param param;
int ret;

/* get the existing scheduling param */
ret = pthread_attr_getschedparam (&tattr, &param);

```

Return Values

`pthread_attr_getschedparam()` returns zero after completing successfully. Any other return value indicates that an error occurred. If the following condition occurs, the function fails and returns the corresponding value.

EINVAL

The value of *param* is NULL or *tattr* is invalid.

Creating a Thread With a Specified Priority

You can set the priority attribute before creating the thread. The child thread is created with the new priority that is specified in the `sched_param` structure (this structure also contains other scheduling information).

It is always a good idea to get the existing parameters, change the priority, create the thread, and then reset the priority.

Example 3-2 shows an example of this.

EXAMPLE 3-2 Creating a Prioritized Thread

```

#include <pthread.h>
#include <sched.h>

pthread_attr_t tattr;
pthread_t tid;
int ret;
int newprio = 20;
sched_param param;

/* initialized with default attributes */
ret = pthread_attr_init (&tattr);

/* safe to get existing scheduling param */
ret = pthread_attr_getschedparam (&tattr, &param);

/* set the priority; others are unchanged */
param.sched_priority = newprio;

/* setting the new scheduling param */
ret = pthread_attr_setschedparam (&tattr, &param);

/* with new priority specified */

```

EXAMPLE 3-2 Creating a Prioritized Thread (Continued)

```
ret = pthread_create (&tid, &tattr, func, arg);
```

Set Stack Size

pthread_attr_setstacksize(3THR)

`pthread_attr_setstacksize(3THR)` sets the thread stack size.

The *stacksize* attribute defines the size of the stack (in bytes) that the system will allocate. The size should not be less than the system-defined minimum stack size. See “About Stacks” on page 65 for more information.

Prototype:

```
int pthread_attr_setstacksize(pthread_attr_t *tattr, size_t size);  
  
#include <pthread.h>  
  
pthread_attr_t tattr;  
size_t size;  
int ret;  
  
size = (PTHREAD_STACK_MIN + 0x4000);  
  
/* setting a new size */  
ret = pthread_attr_setstacksize(&tattr, size);
```

In the example above, *size* contains the number of bytes for the stack that the new thread uses. If *size* is zero, a default size is used. In most cases, a zero value works best.

`PTHREAD_STACK_MIN` is the amount of stack space required to start a thread. This does not take into consideration the threads routine requirements that are needed to execute application code.

Return Values

`pthread_attr_setstacksize()` returns zero after completing successfully. Any other return value indicates that an error occurred. If the following condition occurs, the function fails and returns the corresponding value.

EINVAL

The value of *size* is less than `PTHREAD_STACK_MIN`, or exceeds a system-imposed limit, or *tattr* is not valid.

Get Stack Size

pthread_attr_getstacksize(3THR)

pthread_attr_getstacksize(3THR) returns the stack size set by pthread_attr_setstacksize().

Prototype:

```
int    pthread_attr_getstacksize(pthread_attr_t *tattr, size_t *size);

#include <pthread.h>

pthread_attr_t tattr;
size_t size;
int ret;

/* getting the stack size */
ret = pthread_attr_getstacksize(&tattr, &size);
```

Return Values

pthread_attr_getstacksize() returns zero after completing successfully. Any other return value indicates that an error occurred. If the following condition occurs, the function fails and returns the corresponding value.

EINVAL
tattr is not valid.

About Stacks

Typically, thread stacks begin on page boundaries and any specified size is rounded up to the next page boundary. A page with no access permission is appended to the overflow end of the stack so that most stack overflows result in sending a SIGSEGV signal to the offending thread. Thread stacks allocated by the caller are used as is.

When a stack is specified, the thread should also be created PTHREAD_CREATE_JOINABLE. That stack cannot be freed until the pthread_join(3THR) call for that thread has returned, because the thread's stack cannot be freed until the thread has terminated. The only reliable way to know if such a thread has terminated is through pthread_join(3THR).

Generally, you do not need to allocate stack space for threads. The threads library allocates 1 megabyte (for 32 bit) or 2 megabytes (for 64 bit) of virtual memory for each thread's stack with no swap space reserved. (The library uses the MAP_NORESERVE option of mmap() to make the allocations.)

Each thread stack created by the threads library has a red zone. The library creates the red zone by appending a page to the overflow end of a stack to catch stack overflows. This page is invalid and causes a memory fault if it is accessed. Red zones are appended to all automatically allocated stacks whether the size is specified by the application or the default size is used.

Note – Because runtime stack requirements vary, you should be absolutely certain that the specified stack will satisfy the runtime requirements needed for library calls and dynamic linking.

There are very few occasions when it is appropriate to specify a stack, its size, or both. It is difficult even for an expert to know if the right size was specified. This is because even a program compliant with ABI standards cannot determine its stack size statically. Its size is dependent on the needs of the particular runtime environment in which it executes.

Building Your Own Stack

When you specify the size of a thread stack, be sure to account for the allocations needed by the invoked function and by each function called. The accounting should include calling sequence needs, local variables, and information structures.

Occasionally you want a stack that is a bit different from the default stack. An obvious situation is when the thread needs more than the default stack size. A less obvious situation is when the default stack is too large. You might be creating thousands of threads and not have enough virtual memory to handle the gigabytes of stack space that this many default stacks require.

The limits on the maximum size of a stack are often obvious, but what about the limits on its minimum size? There must be enough stack space to handle all of the stack frames that are pushed onto the stack, along with their local variables, and so on.

You can get the absolute minimum limit on stack size by calling the macro `PTHREAD_STACK_MIN`, which returns the amount of stack space required for a thread that executes a `NULL` procedure. Useful threads need more than this, so be very careful when reducing the stack size.

```
#include <pthread.h>

pthread_attr_t tattr;
pthread_t tid;
int ret;

size_t size = PTHREAD_STACK_MIN + 0x4000;

/* initialized with default attributes */
ret = pthread_attr_init(&tattr);
```

```

/* setting the size of the stack also */
ret = pthread_attr_setstacksize(&tattr, size);

/* only size specified in tattr*/
ret = pthread_create(&tid, &tattr, start_routine, arg);

```

Set Stack Address

pthread_attr_setstackaddr(3THR)

pthread_attr_setstackaddr(3THR) sets the thread stack address.

The *stackaddr* attribute defines the base of the thread's stack. If this is set to non-null (NULL is the default) the system initializes the stack at that address.

Prototype:

```

int pthread_attr_setstackaddr(pthread_attr_t *tattr, void *stackaddr);

#include <pthread.h>

pthread_attr_t tattr;
void *base;
int ret;

base = (void *) malloc(PTHREAD_STACK_MIN + 0x4000);

/* setting a new address */
ret = pthread_attr_setstackaddr(&tattr, base);

```

In the previous example, *base* contains the address for the stack that the new thread uses. If *base* is NULL, then pthread_create(3THR) allocates a stack for the new thread with at least PTHREAD_STACK_MIN bytes.

Return Values

pthread_attr_setstackaddr() returns zero after completing successfully. Any other return value indicates that an error occurred. If the following condition occurs, the function fails and returns the corresponding value.

EINVAL

The value of *base* or *tattr* is incorrect.

This example shows how to create a thread with a custom stack address and size.

```

#include <pthread.h>

pthread_attr_t tattr;

```

```

pthread_t tid;
int ret;
void *stackbase;

/* initialized with default attributes */
ret = pthread_attr_init(&tattr);

/* setting the size of the stack */
ret = pthread_attr_setstacksize(&tattr, size);

/* setting the base address of the stack */
ret = pthread_attr_setstackaddr(&tattr, stackbase);

/* address and size specified */
ret = pthread_create(&tid, &tattr, func, arg);

```

Get Stack Address

pthread_attr_getstackaddr(3THR)

pthread_attr_getstackaddr(3THR) returns the thread stack address set by pthread_attr_setstackaddr().

Prototype:

```

int pthread_attr_getstackaddr(pthread_attr_t *tattr, void * *stackaddr);

#include <pthread.h>

pthread_attr_t tattr;
void *base;
int ret;

/* getting a new address */
ret = pthread_attr_getstackaddr (&tattr, &base);

```

Return Values

pthread_attr_getstackaddr() returns zero after completing successfully. Any other return value indicates that an error occurred. If the following condition occurs, the function fails and returns the corresponding value.

EINVAL

The value or *base* or *tattr* is incorrect.

Programming with Synchronization Objects

This chapter describes the synchronization types available with threads and discusses when and how to use synchronization.

- “Mutual Exclusion Lock Attributes” on page 70
- “Using Mutual Exclusion Locks” on page 86
- “Condition Variable Attributes” on page 98
- “Using Condition Variables” on page 102
- “Semaphores” on page 114
- “Read-Write Lock Attributes” on page 122
- “Set Mutex Attribute’s Protocol” on page 76
- “Synchronization Across Process Boundaries” on page 131
- “Interprocess Locking Without the Threads Library” on page 133
- “Comparing Primitives” on page 133

Synchronization objects are variables in memory that you access just like data. Threads in different processes can communicate with each other through synchronization objects placed in threads-controlled shared memory, even though the threads in different processes are generally invisible to each other.

Synchronization objects can also be placed in files and can have lifetimes beyond that of the creating process.

The available types of synchronization objects are:

- Mutex Locks
- Condition Variables
- Read-Write Locks
- Semaphores

Here are situations that can benefit from the use of synchronization:

- When synchronization is the only way to ensure consistency of shared data.
- When threads in two or more processes can use a single synchronization object jointly. Note that the synchronization object should be initialized by only one of the cooperating processes, because reinitializing a synchronization object sets it to the

unlocked state.

- When synchronization can ensure the safety of mutable data.
- When a process can map a file and have a thread in this process get a record's lock. Once the lock is acquired, any other thread in any process mapping the file that tries to acquire the lock is blocked until the lock is released.
- Even when accessing a single primitive variable, such as an integer. On machines where the integer is not aligned to the bus data width or is larger than the data width, a single memory load can use more than one memory cycle. While this cannot happen on the SPARC[®] Platform Edition architecture, portable programs cannot rely on this.

Note – On 32-bit architectures a long long is not atomic¹ and is read and written as two 32-bit quantities. The types int, char, float, and pointers are atomic on SPARC Platform Edition machines and Intel Architecture machines.

Mutual Exclusion Lock Attributes

Use mutual exclusion locks (mutexes) to serialize thread execution. Mutual exclusion locks synchronize threads, usually by ensuring that only one thread at a time executes a critical section of code. Mutex locks can also preserve single-threaded code.

To change the default mutex attributes, you can declare and initialize an attribute object. Often, the mutex attributes are set in one place at the beginning of the application so they can be located quickly and modified easily. Table 4–1 lists the functions discussed in this section that manipulate mutex attributes.

TABLE 4–1 Mutex Attributes Routines

Operation	Destination Discussion
Initialize a mutex attribute object	"pthread_mutexattr_init(3THR)" on page 72
Destroy a mutex attribute object	"pthread_mutexattr_destroy(3THR)" on page 72
Set the scope of a mutex	"pthread_mutexattr_setpshared(3THR)" on page 73
Get the scope of a mutex	"pthread_mutexattr_getpshared(3THR)" on page 74

¹ An *atomic* operation cannot be divided into smaller operations.

TABLE 4-1 Mutex Attributes Routines (Continued)

Operation	Destination Discussion
Set the mutex type attribute	"pthread_mutexattr_settype(3THR)" on page 75
Get the mutex type attribute	"pthread_mutexattr_gettype(3THR)" on page 76
Set mutex attribute's protocol	"pthread_mutexattr_setprotocol(3THR)" on page 76
Get mutex attribute's protocol	"pthread_mutexattr_getprotocol(3THR)" on page 79
Set mutex attribute's priority ceiling	"pthread_mutexattr_setprioceiling(3THR)" on page 79
Get mutex attribute's priority ceiling	"pthread_mutexattr_getprioceiling(3THR)" on page 80
Set mutex's priority ceiling	"pthread_mutex_setprioceiling(3THR)" on page 82
Get mutex's priority ceiling	"pthread_mutex_getprioceiling(3THR)" on page 83
Set mutex's robust attribute	"pthread_mutexattr_setrobust_np(3THR)" on page 83
Get mutex's robust attribute	"pthread_mutexattr_getrobust_np(3THR)" on page 85

The differences between Solaris threads and POSIX threads, when defining the scope of a mutex, are shown in Table 4-2.

TABLE 4-2 Mutex Scope Comparison

Solaris	POSIX	Definition
USYNC_PROCESS	PTHREAD_PROCESS_SHARED	Use to synchronize threads in this and other processes
USYNC_PROCESS_ROBUST	No POSIX equivalent	Use to <i>robustly</i> synchronize threads between processes
USYNC_THREAD	PTHREAD_PROCESS_PRIVATE	Use to synchronize threads in this process only

Initialize a Mutex Attribute Object

`pthread_mutexattr_init(3THR)`

Use `pthread_mutexattr_init(3THR)` to initialize attributes associated with this object to their default values. Storage for each attribute object is allocated by the threads system during execution.

The default value of the *pshared* attribute when this function is called is `PTHREAD_PROCESS_PRIVATE`, which means that the initialized mutex can be used within a process.

Prototype:

```
int pthread_mutexattr_init(pthread_mutexattr_t *mattr);
```

```
#include <pthread.h>
```

```
pthread_mutexattr_t mattr;  
int ret;
```

```
/* initialize an attribute to default value */
```

```
ret = pthread_mutexattr_init(&mattr);
```

mattr is an opaque type that contains a system-allocated attribute object. The possible values of *mattr*'s scope are `PTHREAD_PROCESS_PRIVATE` (the default) and `PTHREAD_PROCESS_SHARED`.

Before a mutex attribute object can be reinitialized, it must first be destroyed by a call to `pthread_mutexattr_destroy(3THR)`. The `pthread_mutexattr_init()` call results in the allocation of an opaque object. If the object is not destroyed, a memory leak will result.

Return Values

`pthread_mutexattr_init()` returns zero after completing successfully. Any other return value indicates that an error occurred. If either of the following conditions occurs, the function fails and returns the corresponding value.

`ENOMEM`

There is not enough memory to initialize the mutex attributes object.

Destroy a Mutex Attribute Object

`pthread_mutexattr_destroy(3THR)`

`pthread_mutexattr_destroy(3THR)` deallocates the storage space used to maintain the attribute object created by `pthread_mutexattr_init()`.

```

Prototype:
int pthread_mutexattr_destroy(pthread_mutexattr_t *mattr)

#include <pthread.h>

pthread_mutexattr_t mattr;
int ret;

/* destroy an attribute */
ret = pthread_mutexattr_destroy(&mattr);

```

Return Values

`pthread_mutexattr_destroy()` returns zero after completing successfully. Any other return value indicates that an error occurred. If the following condition occurs, the function fails and returns the corresponding value.

EINVAL
The value specified by *mattr* is invalid.

Set the Scope of a Mutex

pthread_mutexattr_setpshared(3THR)

`pthread_mutexattr_setpshared(3THR)` sets the scope of the mutex variable.

The scope of a mutex variable can be either process private (intraprocess) or system wide (interprocess). If the mutex is created with the `pshared` attribute set to the `PTHREAD_PROCESS_SHARED` state, and it exists in shared memory, it can be shared among threads from more than one process. This is equivalent to the `USYNC_PROCESS` flag in `mutex_init()` in the original Solaris threads.

```

Prototype:
int pthread_mutexattr_setpshared(pthread_mutexattr_t *mattr,
    int pshared);

#include <pthread.h>

pthread_mutexattr_t mattr;
int ret;

ret = pthread_mutexattr_init(&mattr);
/*
 * resetting to its default value: private
 */
ret = pthread_mutexattr_setpshared(&mattr,
    PTHREAD_PROCESS_PRIVATE);

```

If the mutex *pshared* attribute is set to `PTHREAD_PROCESS_PRIVATE`, only those threads created by the same process can operate on the mutex.

Return Values

`pthread_mutexattr_setpshared()` returns zero after completing successfully. Any other return value indicates that an error occurred. If the following condition occurs, the function fails and returns the corresponding value.

EINVAL

The value specified by *mattr* is invalid.

Get the Scope of a Mutex

`pthread_mutexattr_getpshared(3THR)`

`pthread_mutexattr_getpshared(3THR)` returns the scope of the mutex variable defined by `pthread_mutexattr_setpshared()`.

Prototype:

```
int pthread_mutexattr_getpshared(pthread_mutexattr_t *mattr,  
int *pshared);
```

```
#include <pthread.h>
```

```
pthread_mutexattr_t mattr;  
int pshared, ret;
```

```
/* get pshared of mutex */
```

```
ret = pthread_mutexattr_getpshared(&mattr, &pshared);
```

Get the current value of *pshared* for the attribute object *mattr*. It is either `PTHREAD_PROCESS_SHARED` or `PTHREAD_PROCESS_PRIVATE`.

Return Values

`pthread_mutexattr_getpshared()` returns zero after completing successfully. Any other return value indicates that an error occurred. If the following condition occurs, the function fails and returns the corresponding value.

EINVAL

The value specified by *mattr* is invalid.

Set the Mutex Type Attribute

pthread_mutexattr_settype(3THR)

```
#include <pthread.h>
```

```
int pthread_mutexattr_settype(pthread_mutexattr_t *attr, int type);
```

pthread_mutexattr_settype(3THR) sets the mutex *type* attribute. The default value of the type attribute is PTHREAD_MUTEX_DEFAULT.

The *type* argument specifies the type of mutex. Valid mutex types include:

PTHREAD_MUTEX_NORMAL

This type of mutex does not detect deadlock. A thread attempting to relock this mutex without first unlocking it will deadlock. Attempting to unlock a mutex locked by a different thread results in undefined behavior. Attempting to unlock an unlocked mutex results in undefined behavior.

PTHREAD_MUTEX_ERRORCHECK

This type of mutex provides error checking. A thread attempting to relock this mutex without first unlocking it will return with an error. A thread attempting to unlock a mutex which another thread has locked will return with an error. A thread attempting to unlock an unlocked mutex will return with an error.

PTHREAD_MUTEX_RECURSIVE

A thread attempting to relock this mutex without first unlocking it will succeed in locking the mutex. The relocking deadlock which can occur with mutexes of type PTHREAD_MUTEX_NORMAL cannot occur with this type of mutex. Multiple locks of this mutex require the same number of unlocks to release the mutex before another thread can acquire the mutex. A thread attempting to unlock a mutex which another thread has locked will return with an error. A thread attempting to unlock an unlocked mutex will return with an error.

PTHREAD_MUTEX_DEFAULT

Attempting to recursively lock a mutex of this type results in undefined behavior. Attempting to unlock a mutex of this type which was not locked by the calling thread results in undefined behavior. Attempting to unlock a mutex of this type which is not locked results in undefined behavior. An implementation is allowed to map this mutex to one of the other mutex types. (For Solaris threads, PTHREAD_PROCESS_DEFAULT is mapped to PTHREAD_PROCESS_NORMAL.)

Return Values

If successful, the pthread_mutexattr_settype function returns zero. Otherwise, an error number is returned to indicate the error.

EINVAL

The value *type* is invalid.

EINVAL

The value specified by *attr* is invalid.

Get the Mutex Type Attribute

pthread_mutexattr_gettype(3THR)

```
#include <pthread.h>
```

```
int pthread_mutexattr_gettype(pthread_mutexattr_t *attr, int *type);
```

`pthread_mutexattr_gettype(3THR)` gets the mutex *type* attribute set by `pthread_mutexattr_settype()`. The default value of the type attribute is `PTHREAD_MUTEX_DEFAULT`.

The *type* argument specifies the type of mutex. Valid mutex types include:

- `PTHREAD_MUTEX_NORMAL`
- `PTHREAD_MUTEX_ERRORCHECK`
- `PTHREAD_MUTEX_RECURSIVE`
- `PTHREAD_MUTEX_DEFAULT`

For a description of each type, see “`pthread_mutexattr_settype(3THR)`” on page 75.

Set Mutex Attribute’s Protocol

pthread_mutexattr_setprotocol(3THR)

`pthread_mutexattr_setprotocol(3THR)` sets the protocol attribute of a mutex attribute object.

```
#include <pthread.h>
```

```
int pthread_mutexattr_setprotocol(pthread_mutexattr_t *attr, int protocol);
```

attr points to a mutex attribute object created by an earlier call to `pthread_mutexattr_init()`.

protocol defines the protocol applied to the mutex attribute object.

The value of *protocol*, defined in `pthread.h`, can be: `PTHREAD_PRIO_NONE`, `PTHREAD_PRIO_INHERIT`, or `PTHREAD_PRIO_PROTECT`.

- `PTHREAD_PRIO_NONE`

A thread's priority and scheduling are not affected by the mutex ownership.

- `PTHREAD_PRIO_INHERIT`

This protocol value affects a thread's (such as `thrd1`) priority and scheduling when higher-priority threads block on one or more mutexes owned by `thrd1` where those mutexes are initialized with `PTHREAD_PRIO_INHERIT`. `thrd1` runs with the higher of its priority or the highest priority of any thread waiting on any of the mutexes owned by `thrd1`.

If `thrd1` blocks on a mutex owned by another thread, `thrd3`, the same priority inheritance effect recursively propagates to `thrd3`.

Use `PTHREAD_PRIO_INHERIT` to avoid priority inversion. Priority inversion occurs when a low-priority thread holds a lock that a higher-priority thread wants. Because the higher-priority thread cannot continue until the lower-priority thread releases the lock, each thread is treated as if it had the inverse of its intended priority.

If the symbol `_POSIX_THREAD_PRIO_INHERIT` is defined, for a mutex initialized with the protocol attribute value `PTHREAD_PRIO_INHERIT`, the following actions occur in the Solaris Operating Environment when the owner of that mutex dies:

Note – The behavior on owner death depends on the value of the *robustness* argument of `pthread_mutexattr_setrobust_np()`.

- The mutex is unlocked.
- The next owner of the mutex acquires it with an error return of `EOWNERDEAD`.
- The next owner of the mutex should try to make the state protected by the mutex consistent—the state might have been left inconsistent when the previous owner died. If the owner is successful in making the state consistent, call `pthread_mutex_init()` for the mutex and unlock the mutex.

Note – If `pthread_mutex_init()` is called on a previously initialized, but not yet destroyed mutex, the mutex is not reinitialized.

- If the owner is unable to make the state consistent, do *not* call `pthread_mutex_init()`, but unlock the mutex. In this event, all waiters will be woken up and all subsequent calls to `pthread_mutex_lock()` will fail to acquire the mutex and return an error code of `ENOTRECOVERABLE`. You can now make the mutex state consistent by calling `pthread_mutex_destroy()` to uninitialized the mutex and calling `pthread_mutex_init()` to reinitialize it.
- If the thread that acquired the lock with `EOWNERDEAD` dies, the next owner acquires the lock with an error code of `EOWNERDEAD`.
- `PTHREAD_PRIO_PROTECT`

This protocol value affects a thread's (such as `thrd2`) priority and scheduling when the thread owns one or more mutexes initialized with `PTHREAD_PRIO_PROTECT`. `thrd2` runs with the higher of its priority or the highest-priority ceiling of all mutexes it owns. Higher-priority threads blocked on any of the mutexes, owned by `thrd2`, have no effect on the scheduling of `thrd2`.

When a thread owns a mutex that is initialized with `PTHREAD_PRIO_INHERIT` or `PTHREAD_PRIO_PROTECT`, and that thread's original priority changes, such as by a call to `sched_setparam()`, the scheduler does not move the thread to the tail of the scheduling queue at its new priority. Similarly, when a thread unlocks a mutex that is initialized with `PTHREAD_PRIO_INHERIT` or `PTHREAD_PRIO_PROTECT`, and that thread's original priority changes, the scheduler does not move the thread to the tail of the scheduling queue at its new priority.

If a thread simultaneously owns several mutexes initialized with a mix of `PTHREAD_PRIO_INHERIT` and `PTHREAD_PRIO_PROTECT`, it executes at the highest priority obtained by either of these protocols.

Return Values

On successful completion, `pthread_mutexattr_setprotocol()` returns 0. Any other return value indicates that an error occurred.

If either of the following conditions occurs, `pthread_mutexattr_setprotocol()` fails and returns the corresponding value.

ENOSYS

Neither of the options `_POSIX_THREAD_PRIO_INHERIT` and `_POSIX_THREAD_PRIO_PROTECT` is defined and the implementation does not support the function.

ENOTSUP

The value specified by *protocol* is an unsupported value.

If either of the following conditions occurs, `pthread_mutexattr_setprotocol()` might fail and return the corresponding value.

EINVAL

The value specified by *attr* or *protocol* is not valid.

EPERM

The caller does not have the privilege to perform the operation.

Get Mutex Attribute's Protocol

`pthread_mutexattr_getprotocol(3THR)`

`pthread_mutexattr_getprotocol(3THR)` gets the protocol attribute of a mutex attribute object.

```
#include <pthread.h>
```

```
int pthread_mutexattr_getprotocol(const pthread_mutexattr_t *attr,  
                                 int *protocol);
```

`attr` points to a mutex attribute object created by an earlier call to `pthread_mutexattr_init()`.

`protocol` contains the protocol attribute: `PTHREAD_PRIO_NONE`, `PTHREAD_PRIO_INHERIT`, or `PTHREAD_PRIO_PROTECT`.

Return Values

On successful completion, `pthread_mutexattr_getprotocol()` returns 0. Any other return value indicates that an error occurred.

If the following condition occurs, `pthread_mutexattr_getprotocol()` fails and returns the corresponding value.

ENOSYS

Neither of the options, `_POSIX_THREAD_PRIO_INHERIT` nor `_POSIX_THREAD_PRIO_PROTECT` is defined and the implementation does not support the function.

If either of the following conditions occurs, `pthread_mutexattr_getprotocol()` might fail and return the corresponding value.

EINVAL

The value specified by `attr` is invalid.

EPERM

The caller does not have the privilege to perform the operation.

Set Mutex Attribute's Priority Ceiling

`pthread_mutexattr_setprioceiling(3THR)`

`pthread_mutexattr_setprioceiling(3THR)` sets the priority ceiling attribute of a mutex attribute object.

```
#include <pthread.h>

int pthread_mutexattr_setprioceiling(pthread_mutexattr_t *attr,
                                     int prioceiling,
                                     int *oldceiling);
```

attr points to a mutex attribute object created by an earlier call to `pthread_mutexattr_init()`.

prioceiling specifies the priority ceiling of initialized mutexes. The ceiling defines the minimum priority level at which the critical section guarded by the mutex is executed. *prioceiling* will be within the maximum range of priorities defined by `SCHED_FIFO`. To avoid priority inversion, *prioceiling* will be set to a priority higher than or equal to the highest priority of all the threads that might lock the particular mutex.

oldceiling contains the old priority ceiling value.

Return Values

On successful completion, `pthread_mutexattr_setprioceiling()` returns 0. Any other return value indicates that an error occurred.

If any of the following conditions occurs, `pthread_mutexattr_setprioceiling()` fails and returns the corresponding value.

ENOSYS

The option `_POSIX_THREAD_PRIO_PROTECT` is not defined and the implementation does not support the function.

If either of the following conditions occurs, `pthread_mutexattr_setprioceiling()` might fail and return the corresponding value.

EINVAL

The value specified by *attr* or *prioceiling* is invalid.

EPERM

The caller does not have the privilege to perform the operation.

Get Mutex Attribute's Priority Ceiling

`pthread_mutexattr_getprioceiling(3THR)`

`pthread_mutexattr_getprioceiling(3THR)` gets the priority ceiling attribute of a mutex attribute object.

```
#include <pthread.h>

int pthread_mutexattr_getprioceiling(const pthread_mutexattr_t *attr,
                                     int *prioceiling);
```

attr designates the attribute object created by an earlier call to `pthread_mutexattr_init()`.

Note – The *attr* mutex attribute object includes the priority ceiling attribute only if the symbol `_POSIX_THREAD_PRIO_PROTECT` is defined.

`pthread_mutexattr_getprioceiling()` returns the priority ceiling of initialized mutexes in *prioceiling*. The ceiling defines the minimum priority level at which the critical section guarded by the mutex is executed. *prioceiling* will be within the maximum range of priorities defined by `SCHED_FIFO`. To avoid priority inversion, *prioceiling* will be set to a priority higher than or equal to the highest priority of all the threads that might lock the particular mutex.

Return Values

On successful completion, `pthread_mutexattr_getprioceiling()` returns 0. Any other return value indicates that an error occurred.

If any of the following conditions occurs, `pthread_mutexattr_getprioceiling()` fails and returns the corresponding value.

ENOSYS

The option `_POSIX_THREAD_PRIO_PROTECT` is not defined and the implementation does not support the function.

If either of the following conditions occurs, `pthread_mutexattr_getprioceiling()` might fail and return the corresponding value.

EINVAL

The value specified by *attr* is invalid.

EPERM

The caller does not have the privilege to perform the operation.

Set Mutex's Priority Ceiling

`pthread_mutex_setprioceiling(3THR)`

`pthread_mutex_setprioceiling(3THR)` sets the priority ceiling of a mutex.

```
#include <pthread.h>
```

```
int pthread_mutex_setprioceiling(pthread_mutex_t *mutex,  
                                int prioceiling,  
                                int *old_ceiling);
```

`pthread_mutex_setprioceiling()` changes the priority ceiling, *prioceiling*, of a mutex, *mutex*. `pthread_mutex_setprioceiling()` locks the mutex if it is unlocked, or blocks until it can successfully lock the mutex, changes the priority ceiling of the mutex and releases the mutex. The process of locking the mutex need not adhere to the priority protect protocol.

If `pthread_mutex_setprioceiling()` succeeds, the previous value of the priority ceiling is returned in *old_ceiling*. If `pthread_mutex_setprioceiling()` fails, the mutex priority ceiling remains unchanged.

Return Values

On successful completion, `pthread_mutex_setprioceiling()` returns 0. Any other return value indicates that an error occurred.

If the following condition occurs, `pthread_mutexatt_setprioceiling()` fails and returns the corresponding value.

ENOSYS

The option `_POSIX_THREAD_PRIO_PROTECT` is not defined and the implementation does not support the function.

If any of the following conditions occurs, `pthread_mutex_setprioceiling()` might fail and return the corresponding value.

EINVAL

The priority requested by *prioceiling* is out of range.

EINVAL

The value specified by *mutex* does not refer to a currently existing mutex.

ENOSYS

The implementation does not support the priority ceiling protocol for mutexes.

EPERM

The caller does not have the privilege to perform the operation.

Get Mutex's Priority Ceiling

`pthread_mutex_getprioceiling(3THR)`

`pthread_mutex_getprioceiling(3THR)` gets the priority ceiling of a mutex.

```
#include <pthread.h>
```

```
int pthread_mutex_getprioceiling(const pthread_mutex_t *mutex,  
                                int *prioceiling);
```

`pthread_mutex_getprioceiling()` returns the priority ceiling, *prioceiling* of a mutex *mutex*.

Return Values

On successful completion, `pthread_mutex_getprioceiling()` returns 0. Any other return value indicates that an error occurred.

If any of the following conditions occurs, `pthread_mutexattr_getprioceiling()` fails and returns the corresponding value.

ENOSYS

The option `_POSIX_THREAD_PRIO_PROTECT` is not defined and the implementation does not support the function.

If any of the following conditions occurs, `pthread_mutex_getprioceiling()` might fail and return the corresponding value.

EINVAL

The value specified by *mutex* does not refer to a currently existing mutex.

ENOSYS

The implementation does not support the priority ceiling protocol for mutexes.

EPERM

The caller does not have the privilege to perform the operation.

Set Mutex's Robust Attribute

`pthread_mutexattr_setrobust_np(3THR)`

`pthread_mutexattr_setrobust_np(3THR)` sets the robust attribute of a mutex attribute object.

```
#include <pthread.h>

int pthread_mutexattr_setrobust_np(pthread_mutexattr_t *attr,
                                   int *robustness);
```

Note – `pthread_mutexattr_setrobust_np()` applies only if the symbol `_POSIX_THREAD_PRIO_INHERIT` is defined.

attr points to the mutex attribute object previously created by a call to `pthread_mutexattr_init()`.

robustness defines the behavior when the owner of the mutex dies. The value of *robustness*, defined in `pthread.h`, is `PTHREAD_MUTEX_ROBUST_NP` or `PTHREAD_MUTEX_STALLED_NP`. The default value is `PTHREAD_MUTEX_STALLED_NP`.

■ `PTHREAD_MUTEX_ROBUST_NP`

When the owner of the mutex dies, all subsequent calls to `pthread_mutex_lock()` are blocked from progress in an unspecified manner.

■ `PTHREAD_MUTEX_STALLED_NP`

When the owner of the mutex dies, the mutex is unlocked. The next owner of this mutex acquires it with an error return of `EOWNERDEAD`.

Note – Your application must check the return code from `pthread_mutex_lock()` for a mutex of this type.

- The new owner of this mutex should make the state protected by the mutex consistent; this state might have been left inconsistent when the previous owner died.
- If the new owner is able to make the state consistent, call `pthread_mutex_consistent_np()` for the mutex, and unlock the mutex.
- If the new owner is *not* able to make the state consistent, do *not* call `pthread_mutex_consistent_np()` for the mutex, but unlock the mutex. All waiters are woken up and all subsequent calls to `pthread_mutex_lock()` fail to acquire the mutex. The return code is `ENOTRECOVERABLE`. The mutex can be made consistent by calling `pthread_mutex_destroy()` to uninitialized the mutex, and calling `pthread_mutex_init()` to reinitialize it.

If the thread that acquire the lock with `EOWNERDEAD` died, the next owner acquires the lock with an `EOWNERDEAD` return code.

Return Values

On successful completion, `pthread_mutexattr_setrobust_np()` returns 0. Any other return value indicates that an error occurred.

If any of the following conditions occurs, `pthread_mutexattr_setrobust_np()` fails and returns the corresponding value.

ENOSYS

The option `_POSIX_THREAD_PRIO__INHERIT` is not defined or the implementation does not support `pthread_mutexattr_setrobust_np()`.

ENOTSUP

The value specified by *robustness* is not supported.

`pthread_mutexattr_setrobust_np()` might fail if:

EINVAL

The value specified by *attr* or *robustness* is invalid.

Get Mutex's Robust Attribute

`pthread_mutexattr_getrobust_np(3THR)`

`pthread_mutexattr_getrobust_np(3THR)` gets the robust attribute of a mutex attribute object.

```
#include <pthread.h>
```

```
int pthread_mutexattr_getrobust_np(const pthread_mutexattr_t *attr,  
                                   int *robustness);
```

Note – `pthread_mutexattr_getrobust_np()` applies only if the symbol `_POSIX_THREAD_PRIO_INHERIT` is defined.

attr points to the mutex attribute object previously created by a call to `pthread_mutexattr_init()`.

robustness is the value of the robust attribute of a mutex attribute object.

Return Values

On successful completion, `pthread_mutexattr_getrobust_np()` returns 0. Any other return value indicates that an error occurred.

If any of the following conditions occurs, `pthread_mutexattr_getrobust_np()` fails and returns the corresponding value.

ENOSYS

The option `_POSIX_THREAD_PRIO_INHERIT` is not defined or the implementation does not support `pthread_mutexattr_getrobust_np()`.

`pthread_mutexattr_getrobust_np()` might fail if:

EINVAL

The value specified by *attr* or *robustness* is invalid.

Using Mutual Exclusion Locks

Table 4-3 lists the functions discussed in this chapter that manipulate mutex locks.

TABLE 4-3 Routines for Mutual Exclusion Locks

Operation	Destination Discussion
Initialize a mutex	" <code>pthread_mutex_init(3THR)</code> " on page 87
Make mutex consistent	" <code>pthread_mutex_consistent_np(3THR)</code> " on page 88
Lock a mutex	" <code>pthread_mutex_lock(3THR)</code> " on page 89
Unlock a mutex	" <code>pthread_mutex_unlock(3THR)</code> " on page 91
Lock with a nonblocking mutex	" <code>pthread_mutex_trylock(3THR)</code> " on page 92
Destroy a mutex	" <code>pthread_mutex_destroy(3THR)</code> " on page 93

The default scheduling policy, `SCHED_OTHER`, does not specify the order in which threads can acquire a lock. When multiple threads are waiting for a mutex, the order of acquisition is undefined. When there is contention, the default behavior is to unblock threads in priority order.

Initialize a Mutex

pthread_mutex_init(3THR)

Use `pthread_mutex_init(3THR)` to initialize the mutex pointed at by `mp` to its default value (`mattr` is `NULL`), or to specify mutex attributes that have already been set with `pthread_mutexattr_init()`. (For Solaris threads, see “`mutex_init(3THR)`” on page 196.)

```
Prototype:
int pthread_mutex_init(pthread_mutex_t *mp,
    const pthread_mutexattr_t *mattr);

#include <pthread.h>

pthread_mutex_t mp = PTHREAD_MUTEX_INITIALIZER;
pthread_mutexattr_t mattr;
int ret;

/* initialize a mutex to its default value */
ret = pthread_mutex_init(&mp, NULL);

/* initialize a mutex */
ret = pthread_mutex_init(&mp, &mattr);
```

When the mutex is initialized, it is in an unlocked state. The mutex can be in memory shared between processes or in memory private to a process.

Note – The mutex memory must be cleared to zero before initialization.

The effect of `mattr` being `NULL` is the same as passing the address of a default mutex attribute object, but without the memory overhead.

Statically defined mutexes can be initialized directly to have default attributes with the macro `PTHREAD_MUTEX_INITIALIZER`.

A mutex lock must not be reinitialized or destroyed while other threads might be using it. Program failure will result if either action is not done correctly. If a mutex is reinitialized or destroyed, the application must be sure the mutex is not currently in use.

Return Values

`pthread_mutex_init()` returns zero after completing successfully. Any other return value indicates that an error occurred. When any of the following conditions occurs, the function fails and returns the corresponding value.

EBUSY

The implementation has detected an attempt to reinitialize the object referenced by *mp* (a previously initialized, but not yet destroyed mutex).

EINVAL

The *matr* attribute value is invalid. The mutex has not been modified.

EFAULT

The address for the mutex pointed at by *mp* is invalid.

Make Mutex Consistent

`pthread_mutex_consistent_np(3THR)`

```
#include <pthread.h>
int pthread_mutex_consistent_np(pthread_mutex_t *mutex);
```

Note – `pthread_mutex_consistent_np()` applies only if the symbol `_POSIX_THREAD_PRIO_INHERIT` is defined and for mutexes that are initialized with the protocol attribute value `PTHREAD_PRIO_INHERIT`.

If the owner of a mutex dies, the mutex can become inconsistent.

`pthread_mutex_consistent_np` makes the mutex object, *mutex*, consistent after the death of its owner.

Call `pthread_mutex_lock()` to acquire the inconsistent mutex. The `EOWNERDEAD` return value indicates an inconsistent mutex.

Call `pthread_mutex_consistent_np()` while holding the mutex acquired by a previous call to `pthread_mutex_lock()`.

Because the critical section protected by the mutex might have been left in an inconsistent state by the dead owner, make the mutex consistent only if you are able to make the critical section protected by the mutex consistent.

Calls to `pthread_mutex_lock()`, `pthread_mutex_unlock()`, and `pthread_mutex_trylock()` for a consistent mutex behave in the normal manner.

The behavior of `pthread_mutex_consistent_np()` for a mutex that is *not* inconsistent, or that is not held, is undefined.

Return Values

`pthread_mutex_consistent_np()` returns zero after completing successfully. Any other return value indicates that an error occurred. When any of the following conditions occurs, the function fails and returns the corresponding value.

`pthread_mutex_consistent_np()` fails if:

ENOSYS

The option `_POSIX_THREAD_PRIO_INHERIT` is not defined or the implementation does not support `pthread_mutex_consistent_np()`.

`pthread_mutex_consistent_np()` might fail if:

EINVAL

The value specified by *mutex* is invalid.

Lock a Mutex

`pthread_mutex_lock(3THR)`

Prototype:

```
int pthread_mutex_lock(pthread_mutex_t *mutex);
```

```
#include <pthread.h>
```

```
pthread_mutex_t mutex;
```

```
int ret;
```

```
ret = pthread_mutex_lock(&mp); /* acquire the mutex */
```

Use `pthread_mutex_lock(3THR)` to lock the mutex pointed to by *mutex*. When `pthread_mutex_lock()` returns, the mutex is locked and the calling thread is the owner. If the mutex is already locked and owned by another thread, the calling thread blocks until the mutex becomes available. (For Solaris threads, see “`mutex_lock(3THR)`” on page 198.)

If the mutex type is `PTHREAD_MUTEX_NORMAL`, deadlock detection is not provided. Attempting to relock the mutex causes deadlock. If a thread attempts to unlock a mutex that it has not locked or a mutex that is unlocked, undefined behavior results.

If the mutex type is `PTHREAD_MUTEX_ERRORCHECK`, then error checking is provided. If a thread attempts to relock a mutex that it has already locked, an error will be returned. If a thread attempts to unlock a mutex that it has not locked or a mutex that is unlocked, an error will be returned.

If the mutex type is `PTHREAD_MUTEX_RECURSIVE`, then the mutex maintains the concept of a lock count. When a thread successfully acquires a mutex for the first time, the lock count is set to one. Every time a thread relocks this mutex, the lock count is

incremented by one. Each time the thread unlocks the mutex, the lock count is decremented by one. When the lock count reaches zero, the mutex becomes available for other threads to acquire. If a thread attempts to unlock a mutex that it has not locked or a mutex which is unlocked, an error will be returned.

If the mutex type is `PTHREAD_MUTEX_DEFAULT`, attempting to recursively lock the mutex results in undefined behavior. Attempting to unlock the mutex if it was not locked by the calling thread results in undefined behavior. Attempting to unlock the mutex if it is not locked results in undefined behavior.

Return Values

`pthread_mutex_lock()` returns zero after completing successfully. Any other return value indicates that an error occurred. When any of the following conditions occurs, the function fails and returns the corresponding value.

`EAGAIN`

The mutex could not be acquired because the maximum number of recursive locks for mutex has been exceeded.

`EDEADLK`

The current thread already owns the mutex.

If the symbol `_POSIX_THREAD_PRIO_INHERIT` is defined, the mutex is initialized with the protocol attribute value `PTHREAD_PRIO_INHERIT`, and the *robustness* argument of `pthread_mutexattr_setrobust_np()` is `PTHREAD_MUTEX_ROBUST_NP` the function fails and returns:

`EOWNERDEAD`

The last owner of this mutex died while holding the mutex. This mutex is now owned by the caller. The caller must attempt to make the state protected by the mutex consistent.

If the caller is able to make the state consistent, call `pthread_mutex_consistent_np()` for the mutex and unlock the mutex. Subsequent calls to `pthread_mutex_lock()` will behave normally.

If the caller is unable to make the state consistent, do not call `pthread_mutex_init()` for the mutex, but unlock the mutex. Subsequent calls to `pthread_mutex_lock()` fail to acquire the mutex and return an `ENOTRECOVERABLE` error code.

If the owner that acquired the lock with `EOWNERDEAD` dies, the next owner acquires the lock with `EOWNERDEAD`.

ENOTRECOVERABLE

The mutex you are trying to acquire is protecting state left irrecoverable by the mutex's previous owner that died while holding the lock. The mutex has not been acquired. This condition can occur when the lock was previously acquired with EOWNERDEAD and the owner was unable to cleanup the state and had unlocked the mutex without making the mutex state consistent.

ENOMEM

The limit on the number of simultaneously held mutexes has been exceeded.

Unlock a Mutex

pthread_mutex_unlock(3THR)

Use `pthread_mutex_unlock(3THR)` to unlock the mutex pointed to by *mutex*. (For Solaris threads, see “`mutex_unlock(3THR)`” on page 198.)

Prototype:

```
int pthread_mutex_unlock(pthread_mutex_t *mutex);
```

```
#include <pthread.h>
```

```
pthread_mutex_t mutex;  
int ret;
```

```
ret = pthread_mutex_unlock(&mutex); /* release the mutex */
```

`pthread_mutex_unlock()` releases the mutex object referenced by *mutex*. The manner in which a mutex is released is dependent upon the mutex's type attribute. If there are threads blocked on the mutex object referenced by *mutex* when `pthread_mutex_unlock()` is called, resulting in the mutex becoming available, the scheduling policy is used to determine which thread shall acquire the mutex. (In the case of `PTHREAD_MUTEX_RECURSIVE` mutexes, the mutex becomes available when the count reaches zero and the calling thread no longer has any locks on this mutex).

Return Values

`pthread_mutex_unlock()` returns zero after completing successfully. Any other return value indicates that an error occurred. When any of the following conditions occurs, the function fails and returns the corresponding value.

EPERM

The current thread does not own the mutex.

Lock With a Nonblocking Mutex

`pthread_mutex_trylock(3THR)`

Use `pthread_mutex_trylock(3THR)` to attempt to lock the mutex pointed to by *mutex*. (For Solaris threads, see “`mutex_trylock(3THR)`” on page 198.)

Prototype:

```
int pthread_mutex_trylock(pthread_mutex_t *mutex);
```

```
#include <pthread.h>
```

```
pthread_mutex_t mutex;
```

```
int ret;
```

```
ret = pthread_mutex_trylock(&mutex); /* try to lock the mutex */
```

`pthread_mutex_trylock()` is a nonblocking version of `pthread_mutex_lock()`. If the mutex object referenced by *mutex* is currently locked (by any thread, including the current thread), the call returns immediately. Otherwise, the mutex is locked and the calling thread is the owner.

Return Values

`pthread_mutex_trylock()` returns zero after completing successfully. Any other return value indicates that an error occurred. When any of the following conditions occurs, the function fails and returns the corresponding value.

EBUSY

The mutex could not be acquired because the mutex pointed to by *mutex* was already locked.

EAGAIN

The mutex could not be acquired because the maximum number of recursive locks for *mutex* has been exceeded.

If the symbol `_POSIX_THREAD_PRIO_INHERIT` is defined, the mutex is initialized with the protocol attribute value `PTHREAD_PRIO_INHERIT`, and the *robustness* argument of `pthread_mutexattr_setrobust_np()` is `PTHREAD_MUTEX_ROBUST_NP` the function fails and returns:

EOWNERDEAD

The last owner of this mutex died while holding the mutex. This mutex is now owned by the caller. The caller must attempt to make the state protected by the mutex consistent.

If the caller is able to make the state consistent, call

`pthread_mutex_consistent_np()` for the mutex and unlock the mutex. Subsequent calls to `pthread_mutex_lock()` will behave normally.

If the caller is unable to make the state consistent, do not call `pthread_mutex_init()` for the mutex, but unlock the mutex. Subsequent calls to `pthread_mutex_trylock()` fail to acquire the mutex and return an `ENOTRECOVERABLE` error code.

If the owner that acquired the lock with `EOWNERDEAD` dies, the next owner acquires the lock with `EOWNERDEAD`.

`ENOTRECOVERABLE`

The mutex you are trying to acquire is protecting state left irrecoverable by the mutex's previous owner that died while holding the lock. The mutex has not been acquired. This condition can occur when the lock was previously acquired with `EOWNERDEAD` and the owner was unable to cleanup the state and had unlocked the mutex without making the mutex state consistent.

`ENOMEM`

The limit on the number of simultaneously held mutexes has been exceeded.

Destroy a Mutex

`pthread_mutex_destroy(3THR)`

Use `pthread_mutex_destroy(3THR)` to destroy any state associated with the mutex pointed to by *mp*. (For Solaris threads, see “`mutex_destroy(3THR)`” on page 197.)

Prototype:

```
int    pthread_mutex_destroy(pthread_mutex_t *mp);

#include <pthread.h>

pthread_mutex_t mp;
int ret;

ret = pthread_mutex_destroy(&mp); /* mutex is destroyed */
```

Note that the space for storing the mutex is not freed.

Return Values

`pthread_mutex_destroy()` returns zero after completing successfully. Any other return value indicates that an error occurred. When any of the following conditions occur, the function fails and returns the corresponding value.

`EINVAL`

The value specified by *mp* does not refer to an initialized mutex object.

Mutex Lock Code Examples

Example 4-1 shows some code fragments with mutex locking.

EXAMPLE 4-1 Mutex Lock Example

```
#include <pthread.h>

pthread_mutex_t count_mutex;
long long count;

void
increment_count()
{
    pthread_mutex_lock(&count_mutex);
    count = count + 1;
    pthread_mutex_unlock(&count_mutex);
}

long long
get_count()
{
    long long c;

    pthread_mutex_lock(&count_mutex);
    c = count;
    pthread_mutex_unlock(&count_mutex);
    return (c);
}
```

The two functions in Example 4-1 use the mutex lock for different purposes. The `increment_count()` function uses the mutex lock simply to ensure an atomic update of the shared variable. The `get_count()` function uses the mutex lock to guarantee that the 64-bit quantity `count` is read atomically. On a 32-bit architecture, a `long long` is really two 32-bit quantities.

Reading an integer value is an atomic operation because integer is the common word size on most machines.

Using Locking Hierarchies

You will occasionally want to access two resources at once. Perhaps you are using one of the resources, and then discover that the other resource is needed as well. There could be a problem if two threads attempt to claim both resources but lock the associated mutexes in different orders. For example, if the two threads lock mutexes 1 and 2 respectively, then a deadlock occurs when each attempts to lock the other mutex. Example 4-2 shows possible deadlock scenarios.

EXAMPLE 4-2 Deadlock

Thread 1	Thread 2
<code>pthread_mutex_lock(&m1);</code>	<code>pthread_mutex_lock(&m2);</code>
<code>/* use resource 1 */</code>	<code>/* use resource 2 */</code>
<code>pthread_mutex_lock(&m2);</code>	<code>pthread_mutex_lock(&m1);</code>
<code>/* use resources 1 and 2 */</code>	<code>/* use resources 1 and 2 */</code>
<code>pthread_mutex_unlock(&m2);</code>	<code>pthread_mutex_unlock(&m1);</code>
<code>pthread_mutex_unlock(&m1);</code>	<code>pthread_mutex_unlock(&m2);</code>

The best way to avoid this problem is to make sure that whenever threads lock multiple mutexes, they do so in the same order. When locks are always taken in a prescribed order, deadlock should not occur. This technique is known as lock hierarchies: order the mutexes by logically assigning numbers to them.

Also, honor the restriction that you cannot take a mutex that is assigned n when you are holding any mutex assigned a number greater than n .

However, this technique cannot always be used—sometimes you must take the mutexes in an order other than prescribed. To prevent deadlock in such a situation, use `pthread_mutex_trylock()`. One thread must release its mutexes when it discovers that deadlock would otherwise be inevitable.

EXAMPLE 4-3 Conditional Locking

Thread 1	Thread 2
<code>pthread_mutex_lock(&m1);</code>	<code>for (; ;)</code>
<code>pthread_mutex_lock(&m2);</code>	<code>{ pthread_mutex_lock(&m2);</code>
<code>/* no processing */</code>	<code>if (pthread_mutex_trylock(&m1) == 0)</code>
<code>pthread_mutex_unlock(&m2);</code>	<code>/* got it! */</code>
<code>pthread_mutex_unlock(&m1);</code>	<code>break;</code>
	<code>/* didn't get it */</code>
	<code>pthread_mutex_unlock(&m2);</code>
	<code>}</code>
	<code>/* get locks; no processing */</code>
	<code>pthread_mutex_unlock(&m1);</code>
	<code>pthread_mutex_unlock(&m2);</code>

In Example 4-3, thread 1 locks mutexes in the prescribed order, but thread 2 takes them out of order. To make certain that there is no deadlock, thread 2 has to take mutex 1 very carefully; if it were to block waiting for the mutex to be released, it is likely to have just entered into a deadlock with thread 1.

To ensure this does not happen, thread 2 calls `pthread_mutex_trylock()`, which takes the mutex if it is available. If it is not, thread 2 returns immediately, reporting failure. At this point, thread 2 must release mutex 2, so that thread 1 can lock it, and then release both mutex 1 and mutex 2.

Nested Locking With a Singly Linked List

Example 4-4 and Example 4-5 show how to take three locks at once, but prevent deadlock by taking the locks in a prescribed order.

EXAMPLE 4-4 Singly Linked List Structure

```
typedef struct node1 {
    int value;
    struct node1 *link;
    pthread_mutex_t lock;
} node1_t;

node1_t ListHead;
```

This example uses a singly linked list structure with each node containing a mutex. To remove a node from the list, first search the list starting at *ListHead* (which itself is never removed) until the desired node is found.

To protect this search from the effects of concurrent deletions, lock each node before any of its contents are accessed. Because all searches start at *ListHead*, there is never a deadlock because the locks are always taken in list order.

When the desired node is found, lock both the node and its predecessor since the change involves both nodes. Because the predecessor's lock is always taken first, you are again protected from deadlock. Example 4-5 shows the C code to remove an item from a singly linked list.

EXAMPLE 4-5 Singly Linked List With Nested Locking

```
node1_t *delete(int value)
{
    node1_t *prev, *current;

    prev = &ListHead;
    pthread_mutex_lock(&prev->lock);
    while ((current = prev->link) != NULL) {
        pthread_mutex_lock(&current->lock);
        if (current->value == value) {
            prev->link = current->link;
```

EXAMPLE 4-5 Singly Linked List With Nested Locking (Continued)

```
        pthread_mutex_unlock(&current->lock);
        pthread_mutex_unlock(&prev->lock);
        current->link = NULL;
        return(current);
    }
    pthread_mutex_unlock(&prev->lock);
    prev = current;
}
pthread_mutex_unlock(&prev->lock);
return(NULL);
}
```

Nested Locking With a Circular Linked List

Example 4-6 modifies the previous list structure by converting it into a circular list. There is no longer a distinguished head node; now a thread might be associated with a particular node and might perform operations on that node and its neighbor. Note that lock hierarchies do not work easily here because the obvious hierarchy (following the links) is circular.

EXAMPLE 4-6 Circular Linked List Structure

```
typedef struct node2 {
    int value;
    struct node2 *link;
    pthread_mutex_t lock;
} node2_t;
```

Here is the C code that acquires the locks on two nodes and performs an operation involving both of them.

EXAMPLE 4-7 Circular Linked List With Nested Locking

```
void Hit Neighbor(node2_t *me) {
    while (1) {
        pthread_mutex_lock(&me->lock);
        if (pthread_mutex_lock(&me->link->lock) != 0) {
            /* failed to get lock */
            pthread_mutex_unlock(&me->lock);
            continue;
        }
        break;
    }
    me->link->value += me->value;
    me->value /=2;
    pthread_mutex_unlock(&me->link->lock);
    pthread_mutex_unlock(&me->lock);
}
```

Condition Variable Attributes

Use condition variables to atomically block threads until a particular condition is true. Always use condition variables together with a mutex lock.

With a condition variable, a thread can atomically block until a condition is satisfied. The condition is tested under the protection of a mutual exclusion lock (mutex).

When the condition is false, a thread usually blocks on a condition variable and atomically releases the mutex waiting for the condition to change. When another thread changes the condition, it can signal the associated condition variable to cause one or more waiting threads to wake up, acquire the mutex again, and reevaluate the condition.

Condition variables can be used to synchronize threads among processes when they are allocated in memory that can be written to and is shared by the cooperating processes.

The scheduling policy determines how blocking threads are awakened. For the default SCHED_OTHER, threads are awakened in priority order.

The attributes for condition variables must be set and initialized before the condition variables can be used. The functions that manipulate condition variable attributes are listed in Table 4-4.

TABLE 4-4 Condition Variable Attributes

Operation	Destination Discussion
Initialize a condition variable attribute	"pthread_condattr_init(3THR)" on page 99
Remove a condition variable attribute	"pthread_condattr_destroy(3THR)" on page 100
Set the scope of a condition variable	"pthread_condattr_setpshared(3THR)" on page 101
Get the scope of a condition variable	"pthread_condattr_getpshared(3THR)" on page 102

The differences between Solaris and POSIX threads, when defining the scope of a condition variable, are shown in Table 4-5.

TABLE 4-5 Condition Variable Scope Comparison

Solaris	POSIX	Definition
USYNC_PROCESS	PTHREAD_PROCESS_SHARED	Use to synchronize threads in this and other processes
USYNC_THREAD	PTHREAD_PROCESS_PRIVATE	Use to synchronize threads in this process only

Initialize a Condition Variable Attribute

pthread_condattr_init(3THR)

Use `pthread_condattr_init(3THR)` to initialize attributes associated with this object to their default values. Storage for each attribute object is allocated by the threads system during execution. The default value of the *pshared* attribute when this function is called is `PTHREAD_PROCESS_PRIVATE`, which means that the initialized condition variable can be used within a process.

Prototype:

```
int pthread_condattr_init(pthread_condattr_t *cattr);
```

```
#include <pthread.h>
pthread_condattr_t cattr;
int ret;
```

```
/* initialize an attribute to default value */
ret = pthread_condattr_init(&cattr);
```

catrr is an opaque data type that contains a system-allocated attribute object. The possible values of *catrr*'s scope are `PTHREAD_PROCESS_PRIVATE` (the default) and `PTHREAD_PROCESS_SHARED`.

Before a condition variable attribute can be reused, it must first be reinitialized by `pthread_condattr_destroy(3THR)`. The `pthread_condattr_init()` call returns a pointer to an opaque object. If the object is not destroyed, a memory leak will result.

Return Values

`pthread_condattr_init()` returns zero after completing successfully. Any other return value indicates that an error occurred. When either of the following conditions occurs, the function fails and returns the corresponding value.

`ENOMEM`

There is not enough memory to initialize the thread attributes object.

`EINVAL`

The value specified by *catrr* is invalid.

Remove a Condition Variable Attribute

`pthread_condattr_destroy(3THR)`

Use `pthread_condattr_destroy(3THR)` to remove storage and render the attribute object invalid.

Prototype:

```
int pthread_condattr_destroy(pthread_condattr_t *catrr);
```

```
#include <pthread.h>
pthread_condattr_t catrr;
int ret;
```

```
/* destroy an attribute */
ret
= pthread_condattr_destroy(&catrr);
```

Return Values

`pthread_condattr_destroy()` returns zero after completing successfully. Any other return value indicates that an error occurred. If the following condition occurs, the function fails and returns the corresponding value.

EINVAL

The value specified by *cattr* is invalid.

Set the Scope of a Condition Variable

pthread_condattr_setpshared(3THR)

`pthread_condattr_setpshared(3THR)` sets the scope of a condition variable to either process private (intraprocess) or system wide (interprocess). If the condition variable is created with the `pshared` attribute set to the `PTHREAD_PROCESS_SHARED` state, and it exists in shared memory, it can be shared among threads from more than one process. This is equivalent to the `USYNC_PROCESS` flag in `mutex_init()` in the original Solaris threads.

If the mutex `pshared` attribute is set to `PTHREAD_PROCESS_PRIVATE` (default value), only those threads created by the same process can operate on the mutex. Using `PTHREAD_PROCESS_PRIVATE` results in the same behavior as with the `USYNC_THREAD` flag in the original Solaris threads `cond_init()` call, which is that of a local condition variable. `PTHREAD_PROCESS_SHARED` is equivalent to a global condition variable.

Prototype:

```
int pthread_condattr_setpshared(pthread_condattr_t *cattr,
                               int pshared);

#include <pthread.h>

pthread_condattr_t cattr;
int ret;

/* all processes */
ret = pthread_condattr_setpshared(&cattr, PTHREAD_PROCESS_SHARED);

/* within a process */
ret = pthread_condattr_setpshared(&cattr, PTHREAD_PROCESS_PRIVATE);
```

Return Values

`pthread_condattr_setpshared()` returns zero after completing successfully. Any other return value indicates that an error occurred. If the following condition occurs, the function fails and returns the corresponding value.

EINVAL

The value of *cattr* is invalid, or the *pshared* value is invalid.

Get the Scope of a Condition Variable

pthread_condattr_getpshared(3THR)

`pthread_condattr_getpshared(3THR)` gets the current value of *pshared* for the attribute object *cattr*. The value is either `PTHREAD_PROCESS_SHARED` or `PTHREAD_PROCESS_PRIVATE`.

```
Prototype:
int pthread_condattr_getpshared(const pthread_condattr_t *cattr,
                               int *pshared);

#include <pthread.h>

pthread_condattr_t cattr;
int pshared;
int ret;

/* get pshared value of condition variable */
ret = pthread_condattr_getpshared(&cattr, &pshared);
```

Return Values

`pthread_condattr_getpshared()` returns zero after completing successfully. Any other return value indicates that an error occurred. When the following condition occurs, the function fails and returns the corresponding value.

EINVAL
The value of *cattr* is invalid.

Using Condition Variables

This section explains using condition variables. Table 4-6 lists the functions that are available.

TABLE 4-6 Condition Variables Functions

Operation	Destination Discussion
Initialize a condition variable	" <code>pthread_cond_init(3THR)</code> " on page 103
Block on a condition variable	" <code>pthread_cond_wait(3THR)</code> " on page 104

TABLE 4-6 Condition Variables Functions (Continued)

Operation	Destination Discussion
Unblock a specific thread	"pthread_cond_signal(3THR)" on page 105
Block until a specified time	"pthread_cond_timedwait(3THR)" on page 107
Block for a specified interval	"pthread_cond_reltimedwait_np(3THR)" on page 108
Unblock all threads	"pthread_cond_broadcast(3THR)" on page 109
Destroy condition variable state	"pthread_cond_destroy(3THR)" on page 110

Initialize a Condition Variable

pthread_cond_init(3THR)

Use `pthread_cond_init(3THR)` to initialize the condition variable pointed at by `cv` to its default value (`cattr` is `NULL`), or to specify condition variable attributes that are already set with `pthread_condattr_init()`. The effect of `cattr` being `NULL` is the same as passing the address of a default condition variable attribute object, but without the memory overhead. (For Solaris threads, see "`cond_init(3THR)`" on page 199.)

Prototype:

```
int pthread_cond_init(pthread_cond_t *cv,
    const pthread_condattr_t *cattr);

#include <pthread.h>

pthread_cond_t cv;
pthread_condattr_t cattr;
int ret;

/* initialize a condition variable to its default value */
ret = pthread_cond_init(&cv, NULL);

/* initialize a condition variable */
ret = pthread_cond_init(&cv, &cattr);
```

Statically defined condition variables can be initialized directly to have default attributes with the macro `PTHREAD_COND_INITIALIZER`. This has the same effect as dynamically allocating `pthread_cond_init()` with null attributes. No error checking is done.

Multiple threads must not simultaneously initialize or reinitialize the same condition variable. If a condition variable is reinitialized or destroyed, the application must be sure the condition variable is not in use.

Return Values

`pthread_cond_init()` returns zero after completing successfully. Any other return value indicates that an error occurred. When any of the following conditions occurs, the function fails and returns the corresponding value.

EINVAL

The value specified by *catrr* is invalid.

EBUSY

The condition variable is being used.

EAGAIN

The necessary resources are not available.

ENOMEM

There is not enough memory to initialize the condition variable.

Block on a Condition Variable

`pthread_cond_wait(3THR)`

Use `pthread_cond_wait(3THR)` to atomically release the mutex pointed to by *mp* and to cause the calling thread to block on the condition variable pointed to by *cv*. (For Solaris threads, see “`cond_wait(3THR)`” on page 200.)

Prototype:

```
int pthread_cond_wait(pthread_cond_t *cv, pthread_mutex_t *mutex);
```

```
#include <pthread.h>
```

```
pthread_cond_t cv;  
pthread_mutex_t mp;  
int ret;
```

```
/* wait on condition variable */  
ret = pthread_cond_wait(&cv, &mp);
```

The blocked thread can be awakened by a `pthread_cond_signal()`, a `pthread_cond_broadcast()`, or when interrupted by delivery of a signal.

Any change in the value of a condition associated with the condition variable cannot be inferred by the return of `pthread_cond_wait()`, and any such condition must be reevaluated.

The `pthread_cond_wait()` routine always returns with the mutex locked and owned by the calling thread, even when returning an error.

This function blocks until the condition is signaled. It atomically releases the associated mutex lock before blocking, and atomically acquires it again before returning.

In typical use, a condition expression is evaluated under the protection of a mutex lock. When the condition expression is false, the thread blocks on the condition variable. The condition variable is then signaled by another thread when it changes the condition value. This causes one or all of the threads waiting on the condition to unblock and to try to acquire the mutex lock again.

Because the condition can change before an awakened thread reacquires the mutex and returns from `pthread_cond_wait()`, and because a waiting thread can be awakened spuriously, the condition that caused the wait must be retested before continuing execution from the point of the `pthread_cond_wait()`. The recommended test method is to write the condition check as a `while()` loop that calls `pthread_cond_wait()`.

```
pthread_mutex_lock();
while(condition_is_false)
    pthread_cond_wait();
pthread_mutex_unlock();
```

No specific order of acquisition is guaranteed when more than one thread blocks on the condition variable.

Note – `pthread_cond_wait()` is a cancellation point. If a cancel is pending and the calling thread has cancellation enabled, the thread terminates and begins executing its cleanup handlers while continuing to hold the lock.

Return Values

`pthread_cond_wait()` returns zero after completing successfully. Any other return value indicates that an error occurred. When the following condition occurs, the function fails and returns the corresponding value.

EINVAL

The value specified by *cv* or *mp* is invalid.

Unblock One Thread

`pthread_cond_signal(3THR)`

Use `pthread_cond_signal(3THR)` to unblock one thread that is blocked on the condition variable pointed to by *cv*. (For Solaris threads, see “`cond_signal(3THR)`” on page 201.)

```

Prototype:
int pthread_cond_signal(pthread_cond_t *cv);

#include <pthread.h>

pthread_cond_t cv;
int ret;

/* one condition variable is signaled */
ret = pthread_cond_signal(&cv);

```

Call `pthread_cond_signal()` under the protection of the same mutex used with the condition variable being signaled. Otherwise, the condition variable could be signaled between the test of the associated condition and blocking in `pthread_cond_wait()`, which can cause an infinite wait.

The scheduling policy determines the order in which blocked threads are awakened. For `SCHED_OTHER`, threads are awakened in priority order.

When no threads are blocked on the condition variable, calling `pthread_cond_signal()` has no effect.

Return Values

`pthread_cond_signal()` returns zero after completing successfully. Any other return value indicates that an error occurred. When the following condition occurs, the function fails and returns the corresponding value.

EINVAL
cv points to an illegal address.

Example 4-8 shows how to use `pthread_cond_wait()` and `pthread_cond_signal()`.

EXAMPLE 4-8 Using `pthread_cond_wait()` and `pthread_cond_signal()`

```

pthread_mutex_t count_lock;
pthread_cond_t count_nonzero;
unsigned count;

decrement_count()
{
    pthread_mutex_lock(&count_lock);
    while (count == 0)
        pthread_cond_wait(&count_nonzero, &count_lock);
    count = count - 1;
    pthread_mutex_unlock(&count_lock);
}

increment_count()
{
    pthread_mutex_lock(&count_lock);

```

EXAMPLE 4-8 Using `pthread_cond_wait()` and `pthread_cond_signal()`
(Continued)

```
    if (count == 0)
        pthread_cond_signal(&count_nonzero);
    count = count + 1;
    pthread_mutex_unlock(&count_lock);
}
```

Block Until a Specified Time

`pthread_cond_timedwait(3THR)`

Prototype:

```
int pthread_cond_timedwait(pthread_cond_t *cv,
    pthread_mutex_t *mp, const struct timespec *abstime);
```

```
#include <pthread.h>
#include <time.h>
```

```
pthread_cond_t cv;
pthread_mutex_t mp;
timestruct_t abstime;
int ret;
```

```
/* wait on condition variable */
ret = pthread_cond_timedwait(&cv, &mp, &abstime);
```

Use `pthread_cond_timedwait(3THR)` as you would use `pthread_cond_wait()`, except that `pthread_cond_timedwait()` does not block past the time of day specified by *abstime*. `pthread_cond_timedwait()` always returns with the mutex locked and owned by the calling thread, even when it is returning an error. (For Solaris threads, see “`cond_timedwait(3THR)`” on page 201.)

The `pthread_cond_timedwait()` function blocks until the condition is signaled or until the time of day, specified by the last argument, has passed.

Note – `pthread_cond_timedwait()` is also a cancellation point.

Return Values

`pthread_cond_timedwait()` returns zero after completing successfully. Any other return value indicates that an error occurred. When either of the following conditions occurs, the function fails and returns the corresponding value.

EINVAL

cv or *abstime* points to an illegal address.

ETIMEDOUT

The time specified by *abstime* has passed.

The timeout is specified as a time of day so that the condition can be retested efficiently without recomputing the value, as shown in Example 4–9.

EXAMPLE 4–9 Timed Condition Wait

```
pthread_timestruc_t to;
pthread_mutex_t m;
pthread_cond_t c;
...
pthread_mutex_lock(&m);
to.tv_sec = time(NULL) + TIMEOUT;
to.tv_nsec = 0;
while (cond == FALSE) {
    err = pthread_cond_timedwait(&c, &m, &to);
    if (err == ETIMEDOUT) {
        /* timeout, do something */
        break;
    }
}
pthread_mutex_unlock(&m);
```

Block For a Specified Interval

pthread_cond_reltimedwait_np(3THR)

Prototype:

```
int pthread_cond_reltimedwait_np(pthread_cond_t *cv,
    pthread_mutex_t *mp,
    const struct timespec *reltime);
```

```
#include <pthread.h>
#include <time.h>
```

```
pthread_cond_t cv;
pthread_mutex_t mp;
timestruc_t reltime;
int ret;
```

```
/* wait on condition variable */
ret = pthread_cond_reltimedwait_np(&cv, &mp, &reltime);
```

Use `pthread_cond_reltimedwait_np(3THR)` as you would use `pthread_cond_timedwait()`, except that `pthread_cond_reltimedwait_np()` takes a relative time interval rather than an absolute future time of day as the value of its last argument. `pthread_cond_reltimedwait_np()` always returns with the mutex locked and owned by the calling thread, even when it is returning an error. (For

Solaris threads, see `cond_reltimedwait(3THR)`.) The `pthread_cond_reltimedwait_np()` function blocks until the condition is signaled or until the time interval, specified by the last argument, has elapsed.

Note – `pthread_cond_reltimedwait_np()` is also a cancellation point.

Return Values

`pthread_cond_reltimedwait_np()` returns zero after completing successfully. Any other return value indicates that an error occurred. When either of the following conditions occurs, the function fails and returns the corresponding value.

EINVAL

cv or *reltime* points to an illegal address.

ETIMEDOUT

The time interval specified by *reltime* has passed.

Unblock All Threads

`pthread_cond_broadcast(3THR)`

Prototype:

```
int pthread_cond_broadcast(pthread_cond_t *cv);
```

```
#include <pthread.h>
```

```
pthread_cond_t cv;  
int ret;
```

```
/* all condition variables are signaled */  
ret = pthread_cond_broadcast(&cv);
```

Use `pthread_cond_broadcast(3THR)` to unblock all threads that are blocked on the condition variable pointed to by *cv*, specified by `pthread_cond_wait()`. When no threads are blocked on the condition variable, `pthread_cond_broadcast()` has no effect. (For Solaris threads, see “`cond_broadcast(3THR)`” on page 202.)

Return Values

`pthread_cond_broadcast()` returns zero after completing successfully. Any other return value indicates that an error occurred. When the following condition occurs, the function fails and returns the corresponding value.

EINVAL

cv points to an illegal address.

Condition Variable Broadcast Example

Since `pthread_cond_broadcast()` causes all threads blocked on the condition to contend again for the mutex lock, use it with care. For example, use `pthread_cond_broadcast()` to allow threads to contend for varying resource amounts when resources are freed, as shown in Example 4–10.

EXAMPLE 4–10 Condition Variable Broadcast

```
pthread_mutex_t rsrc_lock;
pthread_cond_t rsrc_add;
unsigned int resources;

get_resources(int amount)
{
    pthread_mutex_lock(&rsrc_lock);
    while (resources < amount) {
        pthread_cond_wait(&rsrc_add, &rsrc_lock);
    }
    resources -= amount;
    pthread_mutex_unlock(&rsrc_lock);
}

add_resources(int amount)
{
    pthread_mutex_lock(&rsrc_lock);
    resources += amount;
    pthread_cond_broadcast(&rsrc_add);
    pthread_mutex_unlock(&rsrc_lock);
}
```

Note that in `add_resources()` it does not matter whether *resources* is updated first or if `pthread_cond_broadcast()` is called first inside the mutex lock.

Call `pthread_cond_broadcast()` under the protection of the same mutex that is used with the condition variable being signaled. Otherwise, the condition variable could be signaled between the test of the associated condition and blocking in `pthread_cond_wait()`, which can cause an infinite wait.

Destroy Condition Variable State

`pthread_cond_destroy(3THR)`

Use `pthread_cond_destroy(3THR)` to destroy any state associated with the condition variable pointed to by *cv*. (For Solaris threads, see “`cond_destroy(3THR)`” on page 200.)

```

Prototype:
int pthread_cond_destroy(pthread_cond_t *cv);

#include <pthread.h>

pthread_cond_t cv;
int ret;

/* Condition variable is destroyed */
ret = pthread_cond_destroy(&cv);

```

Note that the space for storing the condition variable is not freed.

Return Values

`pthread_cond_destroy()` returns zero after completing successfully. Any other return value indicates that an error occurred. When any of the following conditions occur, the function fails and returns the corresponding value.

EINVAL

The value specified by *cv* is invalid.

The Lost Wake-Up Problem

Calling `pthread_cond_signal()` or `pthread_cond_broadcast()` when the thread does not hold the mutex lock associated with the condition can lead to *lost wake-up* bugs.

A lost wake-up occurs when:

- A thread calls `pthread_cond_signal()` or `pthread_cond_broadcast()`.
- *And* another thread is between the test of the condition and the call to `pthread_cond_wait()`.
- *And* no threads are waiting.

The signal has no effect, and therefore is lost.

The Producer/Consumer Problem

This problem is one of the small collection of standard, well-known problems in concurrent programming: a finite-size buffer and two classes of threads, producers and consumers, put items into the buffer (producers) and take items out of the buffer (consumers).

A producer must wait until the buffer has space before it can put something in, and a consumer must wait until something is in the buffer before it can take something out.

A condition variable represents a queue of threads waiting for some condition to be signaled.

Example 4–11 has two such queues, one (*less*) for producers waiting for a slot in the buffer, and the other (*more*) for consumers waiting for a buffer slot containing information. The example also has a mutex, as the data structure describing the buffer must be accessed by only one thread at a time.

EXAMPLE 4–11 The Producer/Consumer Problem and Condition Variables

```
typedef struct {
    char buf[BFSIZE];
    int occupied;
    int nextin;
    int nextout;
    pthread_mutex_t mutex;
    pthread_cond_t more;
    pthread_cond_t less;
} buffer_t;

buffer_t buffer;
```

As Example 4–12 shows, the producer thread acquires the mutex protecting the buffer data structure and then makes certain that space is available for the item being produced. If not, it calls `pthread_cond_wait()`, which causes it to join the queue of threads waiting for the condition *less*, representing *there is room in the buffer*, to be signaled.

At the same time, as part of the call to `pthread_cond_wait()`, the thread releases its lock on the mutex. The waiting producer threads depend on consumer threads to signal when the condition is true (as shown in Example 4–12). When the condition is signaled, the first thread waiting on *less* is awakened. However, before the thread can return from `pthread_cond_wait()`, it must acquire the lock on the mutex again.

This ensures that it again has mutually exclusive access to the buffer data structure. The thread then must check that there really is room available in the buffer; if so, it puts its item into the next available slot.

At the same time, consumer threads might be waiting for items to appear in the buffer. These threads are waiting on the condition variable *more*. A producer thread, having just deposited something in the buffer, calls `pthread_cond_signal()` to wake up the next waiting consumer. (If there are no waiting consumers, this call has no effect.)

Finally, the producer thread unlocks the mutex, allowing other threads to operate on the buffer data structure.

EXAMPLE 4–12 The Producer/Consumer Problem—the Producer

```
void producer(buffer_t *b, char item)
{
    pthread_mutex_lock(&b->mutex);
```

EXAMPLE 4-12 The Producer/Consumer Problem—the Producer (Continued)

```
while (b->occupied >= BSIZE)
    pthread_cond_wait(&b->less, &b->mutex);

assert(b->occupied < BSIZE);

b->buf[b->nextin++] = item;

b->nextin %= BSIZE;
b->occupied++;

/* now: either b->occupied < BSIZE and b->nextin is the index
   of the next empty slot in the buffer, or
   b->occupied == BSIZE and b->nextin is the index of the
   next (occupied) slot that will be emptied by a consumer
   (such as b->nextin == b->nextout) */

pthread_cond_signal(&b->more);

pthread_mutex_unlock(&b->mutex);
}
```

Note the use of the `assert()` statement; unless the code is compiled with `NDEBUG` defined, `assert()` does nothing when its argument evaluates to true (that is, nonzero), but causes the program to abort if the argument evaluates to false (zero). Such assertions are especially useful in multithreaded programs—they immediately point out runtime problems if they fail, and they have the additional effect of being useful comments.

The comment that begins `/* now: either b->occupied ...` could better be expressed as an assertion, but it is too complicated as a Boolean-valued expression and so is given in English.

Both the assertion and the comments are examples of invariants. These are logical statements that should not be falsified by the execution of the program, except during brief moments when a thread is modifying some of the program variables mentioned in the invariant. (An assertion, of course, should be true whenever any thread executes it.)

Using invariants is an extremely useful technique. Even if they are not stated in the program text, think in terms of invariants when you analyze a program.

The invariant in the producer code that is expressed as a comment is always true whenever a thread is in the part of the code where the comment appears. If you move this comment to just after the `mutex_unlock()`, this does not necessarily remain true. If you move this comment to just after the `assert()`, this is still true.

The point is that this invariant expresses a property that is true at all times, except when either a producer or a consumer is changing the state of the buffer. While a thread is operating on the buffer (under the protection of a mutex), it might temporarily falsify the invariant. However, once the thread is finished, the invariant should be true again.

Example 4–13 shows the code for the consumer. Its flow is symmetric with that of the producer.

EXAMPLE 4–13 The Producer/Consumer Problem—the Consumer

```
char consumer(buffer_t *b)
{
    char item;
    pthread_mutex_lock(&b->mutex);
    while(b->occupied <= 0)
        pthread_cond_wait(&b->more, &b->mutex);

    assert(b->occupied > 0);

    item = b->buf[b->nextout++];
    b->nextout %= BSIZE;
    b->occupied--;

    /* now: either b->occupied > 0 and b->nextout is the index
       of the next occupied slot in the buffer, or
       b->occupied == 0 and b->nextout is the index of the next
       (empty) slot that will be filled by a producer (such as
       b->nextout == b->nextin) */

    pthread_cond_signal(&b->less);
    pthread_mutex_unlock(&b->mutex);

    return(item);
}
```

Semaphores

Semaphores are a programming construct designed by E. W. Dijkstra in the late 1960s. Dijkstra's model was the operation of railroads: consider a stretch of railroad in which there is a single track over which only one train at a time is allowed.

Guarding this track is a semaphore. A train must wait before entering the single track until the semaphore is in a state that permits travel. When the train enters the track, the semaphore changes state to prevent other trains from entering the track. A train that is leaving this section of track must again change the state of the semaphore to allow another train to enter.

In the computer version, a semaphore appears to be a simple integer. A thread waits for permission to proceed and then signals that it has proceeded by performing a P operation on the semaphore.

The semantics of the operation are such that the thread must wait until the semaphore's value is positive, then change the semaphore's value by subtracting one from it. When it is finished, the thread performs a V operation, which changes the semaphore's value by adding one to it. It is crucial that these operations take place atomically—they cannot be subdivided into pieces between which other actions on the semaphore can take place. In the P operation, the semaphore's value must be positive just before it is decremented (resulting in a value that is guaranteed to be nonnegative and one less than what it was before it was decremented).

In both P and V operations, the arithmetic must take place without interference. If two V operations are performed simultaneously on the same semaphore, the net effect should be that the semaphore's new value is two greater than it was.

The mnemonic significance of P and V is unclear to most of the world, as Dijkstra is Dutch. However, in the interest of true scholarship: P stands for *prolagen*, a made-up word derived from *proberen te verlagen*, which means *try to decrease*. V stands for *verhogen*, which means *increase*. This is discussed in one of Dijkstra's technical notes, EWD 74.

`sem_wait(3RT)` and `sem_post(3RT)` correspond to Dijkstra's P and V operations. `sem_trywait(3RT)` is a conditional form of the P operation: if the calling thread cannot decrement the value of the semaphore without waiting, the call returns immediately with a nonzero value.

There are two basic sorts of semaphores: binary semaphores, which never take on values other than zero or one, and counting semaphores, which can take on arbitrary nonnegative values. A binary semaphore is logically just like a mutex.

However, although it is not enforced, mutexes should be unlocked only by the thread holding the lock. There is no notion of "the thread holding the semaphore," so any thread can perform a V (or `sem_post(3RT)`) operation.

Counting semaphores are about as powerful as conditional variables (used in conjunction with mutexes). In many cases, the code might be simpler when it is implemented with counting semaphores rather than with condition variables (as shown in the next few examples).

However, when a mutex is used with condition variables, there is an implied bracketing—it is clear which part of the program is being protected. This is not necessarily the case for a semaphore, which might be called the *go to* of concurrent programming—it is powerful but too easy to use in an unstructured, indeterminate way.

Counting Semaphores

Conceptually, a semaphore is a nonnegative integer count. Semaphores are typically used to coordinate access to resources, with the semaphore count initialized to the number of free resources. Threads then atomically increment the count when resources are added and atomically decrement the count when resources are removed.

When the semaphore count becomes zero, indicating that no more resources are present, threads trying to decrement the semaphore block wait until the count becomes greater than zero.

TABLE 4-7 Routines for Semaphores

Operation	Destination Discussion
Initialize a semaphore	"sem_init(3RT)" on page 116
Increment a semaphore	"sem_post(3RT)" on page 118
Block on a semaphore count	"sem_wait(3RT)" on page 119
Decrement a semaphore count	"sem_trywait(3RT)" on page 119
Destroy the semaphore state	"sem_destroy(3RT)" on page 120

Because semaphores need not be acquired and released by the same thread, they can be used for asynchronous event notification (such as in signal handlers). And, because semaphores contain state, they can be used asynchronously without acquiring a mutex lock as is required by condition variables. However, semaphores are not as efficient as mutex locks.

By default, there is no defined order of unblocking if multiple threads are waiting for a semaphore.

Semaphores must be initialized before use, but they do not have attributes.

Initialize a Semaphore

sem_init(3RT)

Prototype:

```
int sem_init(sem_t *sem, int pshared, unsigned int value);
```

```
#include <semaphore.h>
```

```
sem_t sem;  
int pshared;  
int ret;
```

```

int value;

/* initialize a private semaphore */
pshared = 0;
value = 1;
ret = sem_init(&sem, pshared, value);

```

Use `sem_init(3THR)` to initialize the semaphore variable pointed to by `sem` to `value` amount. If the value of `pshared` is zero, then the semaphore cannot be shared between processes. If the value of `pshared` is nonzero, then the semaphore can be shared between processes. (For Solaris threads, see “`sem_init(3THR)`” on page 202.)

Multiple threads must not initialize the same semaphore.

A semaphore must not be reinitialized while other threads might be using it.

Return Values

`sem_init()` returns zero after completing successfully. Any other return value indicates that an error occurred. When any of the following conditions occurs, the function fails and returns the corresponding value.

EINVAL

The value argument exceeds `SEM_VALUE_MAX`.

ENOSPC

A resource required to initialize the semaphore has been exhausted. The limit on semaphores `SEM_NSEMS_MAX` has been reached.

EPERM

The process lacks the appropriate privileges to initialize the semaphore.

Initializing Semaphores With Intraprocess Scope

When `pshared` is 0, the semaphore can be used by all the threads in this process only.

```

#include <semaphore.h>

sem_t sem;
int ret;
int count = 4;

/* to be used within this process only */
ret = sem_init(&sem, 0, count);

```

Initializing Semaphores With Interprocess Scope

When `pshared` is nonzero, the semaphore can be shared by other processes.

```

#include <semaphore.h>

sem_t sem;
int ret;
int count = 4;

/* to be shared among processes */
ret = sem_init(&sem, 1, count);

```

Named Semaphores

The functions `sem_open(3RT)`, `sem_getvalue(3RT)`, `sem_close(3RT)`, and `sem_unlink(3RT)` are available to open, retrieve, close, and remove named semaphores. Using `sem_open()`, you can create a semaphore that has a name defined in the file system name space.

Named semaphores are like process shared semaphores, except that they are referenced with a pathname rather than a *pshared* value.

For more information about named semaphores, see `sem_open(3RT)`, `sem_getvalue(3RT)`, `sem_close(3RT)`, and `sem_unlink(3RT)`.

Increment a Semaphore

`sem_post(3RT)`

```

Prototype:
int    sem_post(sem_t *sem);

#include <semaphore.h>

sem_t sem;
int ret;

ret = sem_post(&sem); /* semaphore is posted */

```

Use `sema_post(3THR)` to atomically increment the semaphore pointed to by *sem*. When any threads are blocked on the semaphore, one of them is unblocked. (For Solaris threads, see “`sema_post(3THR)`” on page 203.)

Return Values

`sem_post()` returns zero after completing successfully. Any other return value indicates that an error occurred. When the following condition occurs, the function fails and returns the corresponding value.

EINVAL
 sem points to an illegal address.

Block on a Semaphore Count

`sem_wait(3RT)`

```
Prototype:  
int    sem_wait(sem_t *sem);  
  
#include <semaphore.h>  
  
sem_t sem;  
int ret;  
  
ret = sem_wait(&sem); /* wait for semaphore */
```

Use `sem_wait(3THR)` to block the calling thread until the count in the semaphore pointed to by *sem* becomes greater than zero, then atomically decrement it.

Return Values

`sem_wait()` returns zero after completing successfully. Any other return value indicates that an error occurred. When any of the following conditions occurs, the function fails and returns the corresponding value.

EINVAL
 sem points to an illegal address.

EINTR
 A signal interrupted this function.

Decrement a Semaphore Count

`sem_trywait(3RT)`

```
Prototype:  
int    sem_trywait(sem_t *sem);  
  
#include <semaphore.h>  
  
sem_t sem;
```

```
int ret;

ret = sem_trywait(&sem); /* try to wait for semaphore*/
```

Use `sem_trywait(3RT)` to try to atomically decrement the count in the semaphore pointed to by `sem` when the count is greater than zero. This function is a nonblocking version of `sem_wait()`; that is it returns immediately if unsuccessful.

Return Values

`sem_trywait()` returns zero after completing successfully. Any other return value indicates that an error occurred. When any of the following conditions occurs, the function fails and returns the corresponding value.

EINVAL
`sem` points to an illegal address.

EINTR
A signal interrupted this function.

EAGAIN
The semaphore was already locked, so it cannot be immediately locked by the `sem_trywait()` operation.

Destroy the Semaphore State

`sem_destroy(3RT)`

```
Prototype:
int    sem_destroy(sem_t *sem);

#include <semaphore.h>

sem_t sem;
int ret;

ret = sem_destroy(&sem); /* the semaphore is destroyed */
```

Use `sem_destroy(3RT)` to destroy any state associated with the semaphore pointed to by `sem`. The space for storing the semaphore is not freed. (For Solaris threads, see “`sem_destroy(3THR)`” on page 204.)

Return Values

`sem_destroy()` returns zero after completing successfully. Any other return value indicates that an error occurred. When the following condition occurs, the function fails and returns the corresponding value.

EINVAL

sem points to an illegal address.

The Producer/Consumer Problem, Using Semaphores

The data structure in Example 4-14 is similar to that used for the condition variables example (see Example 4-11). Two semaphores represent the number of full and empty buffers and ensure that producers wait until there are empty buffers and that consumers wait until there are full buffers.

EXAMPLE 4-14 The Producer/Consumer Problem With Semaphores

```
typedef struct {
    char buf[BFSIZE];
    sem_t occupied;
    sem_t empty;
    int nextin;
    int nextout;
    sem_t pmut;
    sem_t cmut;
} buffer_t;

buffer_t buffer;

sem_init(&buffer.occupied, 0, 0);
sem_init(&buffer.empty, 0, BFSIZE);
sem_init(&buffer.pmut, 0, 1);
sem_init(&buffer.cmut, 0, 1);
buffer.nextin = buffer.nextout = 0;
```

Another pair of (binary) semaphores plays the same role as mutexes, controlling access to the buffer when there are multiple producers and multiple empty buffer slots, and when there are multiple consumers and multiple full buffer slots. Mutexes would work better here, but would not provide as good an example of semaphore use.

EXAMPLE 4-15 The Producer/Consumer Problem—the Producer

```
void producer(buffer_t *b, char item) {
    sem_wait(&b->empty);
    sem_wait(&b->pmut);

    b->buf[b->nextin] = item;
    b->nextin++;
    b->nextin %= BFSIZE;

    sem_post(&b->pmut);
    sem_post(&b->occupied);
}
```

EXAMPLE 4-16 The Producer/Consumer Problem—the Consumer

```
char consumer(buffer_t *b) {
    char item;

    sem_wait(&b->occupied);

    sem_wait(&b->cmut);

    item = b->buf[b->nextout];
    b->nextout++;
    b->nextout %= BSIZE;

    sem_post(&b->cmut);

    sem_post(&b->empty);

    return(item);
}
```

Read-Write Lock Attributes

Read-write locks permit concurrent reads and exclusive writes to a protected shared resource. The read-write lock is a single entity that can be locked in *read* or *write* mode. To modify a resource, a thread must first acquire the exclusive write lock. An exclusive write lock is not permitted until all read locks have been released.

Database access can be synchronized with a read-write lock. Read-write locks support concurrent reads of database records because the read operation does not change the record's information. When the database is to be updated, the write operation must acquire an exclusive write lock.

To change the default read-write lock attributes, you can declare and initialize an attribute object. Often, the read-write lock attributes are set up in one place at the beginning of the application so they can be located quickly and modified easily. The following table lists the functions discussed in this section that manipulate read-write lock attributes.

See “Similar Synchronization Functions—Read-Write Locks” on page 181 for the Solaris threads implementation of read-write locks.

TABLE 4-8 Routines for Read-Write Lock Attributes

Operation	Destination Discussion
Initialize a read-write lock attribute	“pthread_rwlockattr_init(3THR)” on page 123
Destroy a read-write lock attribute	“pthread_rwlockattr_destroy(3THR)” on page 123
Set a read-write lock attribute	“pthread_rwlockattr_setpshared(3THR)” on page 124
Get a read-write lock attribute	“pthread_rwlockattr_getpshared(3THR)” on page 125

Initialize a Read-Write Lock Attribute

pthread_rwlockattr_init(3THR)

```
#include <pthread.h>
```

```
int pthread_rwlockattr_init(pthread_rwlockattr_t *attr);
```

`pthread_rwlockattr_init(3THR)` initializes a read-write lock attributes object *attr* with the default value for all of the attributes defined by the implementation.

Results are undefined if `pthread_rwlockattr_init` is called specifying an already initialized read-write lock attributes object. After a read-write lock attributes object has been used to initialize one or more read-write locks, any function affecting the attributes object (including destruction) does not affect any previously initialized read-write locks.

Return Values

If successful, `pthread_rwlockattr_init()` returns zero. Otherwise, an error number is returned to indicate the error.

ENOMEM

Insufficient memory exists to initialize the rwlock attributes object.

Destroy a Read-Write Lock Attribute

pthread_rwlockattr_destroy(3THR)

```
#include <pthread.h>
```

```
int pthread_rwlockattr_destroy(pthread_rwlockattr_t *attr);
```

`pthread_rwlockattr_destroy(3THR)` destroys a read-write lock attributes object. The effect of subsequent use of the object is undefined until the object is re-initialized by another call to `pthread_rwlockattr_init()`. An implementation can cause `pthread_rwlockattr_destroy()` to set the object referenced by *attr* to an invalid value.

Return Values

If successful, `pthread_rwlockattr_destroy()` returns zero. Otherwise, an error number is returned to indicate the error.

EINVAL

The value specified by *attr* is invalid.

Set a Read-Write Lock Attribute

`pthread_rwlockattr_setpshared(3THR)`

```
#include <pthread.h>
```

```
int pthread_rwlockattr_setpshared(pthread_rwlockattr_t *attr,  
                                int pshared);
```

`pthread_rwlockattr_setpshared(3THR)` sets the process-shared read-write lock attribute.

PTHREAD_PROCESS_SHARED

Permits a read-write lock to be operated on by any thread that has access to the memory where the read-write lock is allocated, even if the read-write lock is allocated in memory that is shared by multiple processes.

PTHREAD_PROCESS_PRIVATE

The read-write lock will only be operated upon by threads created within the same process as the thread that initialized the read-write lock; if threads of differing processes attempt to operate on such a read-write lock, the behavior is undefined. The default value of the process-shared attribute is `PTHREAD_PROCESS_PRIVATE`.

Return Value

If successful, `pthread_rwlockattr_setpshared()` returns zero. Otherwise, an error number is returned to indicate the error.

EINVAL

The value specified by *attr* or *pshared* is invalid.

Get a Read-Write Lock Attribute

pthread_rwlockattr_getpshared(3THR)

```
#include <pthread.h>

int pthread_rwlockattr_getpshared(const pthread_rwlockattr_t *attr,
                                  int *pshared);
```

`pthread_rwlockattr_getpshared(3THR)` gets the process-shared read-write lock attribute.

`pthread_rwlockattr_getpshared()` obtains the value of the process-shared attribute from the initialized attributes object referenced by *attr*.

Return Value

If successful, `pthread_rwlockattr_getpshared()` returns zero. Otherwise, an error number is returned to indicate the error.

EINVAL

The value specified by *attr* or *pshared* is invalid.

Using Read-Write Locks

After the attributes for a read-write lock are configured, you initialize the read-write lock itself. The following functions are used to initialize or destroy, lock or unlock, or try to lock a read-write lock. The following table lists the functions discussed in this section that manipulate read-write locks.

TABLE 4-9 Routines that Manipulate Read-Write Locks

Operation	Destination Discussion
Initialize a read-write lock	“pthread_rwlock_init(3THR)” on page 126
Read lock on read-write lock	“pthread_rwlock_rdlock(3THR)” on page 127
Read lock with a nonblocking read-write lock	“pthread_rwlock_tryrdlock(3THR)” on page 128
Write lock on read-write lock	“pthread_rwlock_wrlock(3THR)” on page 128

TABLE 4-9 Routines that Manipulate Read-Write Locks (Continued)

Operation	Destination Discussion
Write lock with a nonblocking read-write lock	"pthread_rwlock_trywrlock(3THR)" on page 129
Unlock a read-write lock	"pthread_rwlock_unlock(3THR)" on page 129
Destroy a read-write lock	"pthread_rwlock_destroy(3THR)" on page 130

Initialize a Read-Write Lock

pthread_rwlock_init(3THR)

```
#include <pthread.h>

int pthread_rwlock_init(pthread_rwlock_t *rwlock,
                       const pthread_rwlockattr_t *attr);

pthread_rwlock_t  rwlock = PTHREAD_RWLOCK_INITIALIZER;
```

Use `pthread_rwlock_init(3THR)` to initialize the read-write lock referenced by `rwlock` with the attributes referenced by `attr`. If `attr` is `NULL`, the default read-write lock attributes are used; the effect is the same as passing the address of a default read-write lock attributes object. Once initialized, the lock can be used any number of times without being re-initialized. On successful initialization, the state of the read-write lock becomes initialized and unlocked. Results are undefined if `pthread_rwlock_init()` is called specifying an already initialized read-write lock. Results are undefined if a read-write lock is used without first being initialized. (For Solaris threads, see "rwlock_init(3THR)" on page 182.)

In cases where default read-write lock attributes are appropriate, the macro `PTHREAD_RWLOCK_INITIALIZER` can be used to initialize read-write locks that are statically allocated. The effect is equivalent to dynamic initialization by a call to `pthread_rwlock_init()` with the parameter `attr` specified as `NULL`, except that no error checks are performed.

Return Value

If successful, `pthread_rwlock_init()` returns zero. Otherwise, an error number is returned to indicate the error.

If `pthread_rwlock_init()` fails, `rwlock` is not initialized and the contents of `rwlock` are undefined.

EINVAL

The value specified by `attr` or `rwlock` is invalid.

Read Lock on Read-Write Lock

`pthread_rwlock_rdlock(3THR)`

```
#include <pthread.h>
```

```
int pthread_rwlock_rdlock(pthread_rwlock_t *rwlock );
```

`pthread_rwlock_rdlock(3THR)` applies a read lock to the read-write lock referenced by *rwlock*. The calling thread acquires the read lock if a writer does not hold the lock and there are no writers blocked on the lock. It is unspecified whether the calling thread acquires the lock when a writer does not hold the lock and there are writers waiting for the lock. If a writer holds the lock, the calling thread will not acquire the read lock. If the read lock is not acquired, the calling thread blocks (that is, it does not return from the `pthread_rwlock_rdlock()`) until it can acquire the lock. Results are undefined if the calling thread holds a write lock on *rwlock* at the time the call is made.

Implementations are allowed to favor writers over readers to avoid writer starvation. (For instance, the Solaris threads implementation favors writers over readers. See “`rw_rdlock(3THR)`” on page 183.)

A thread can hold multiple concurrent read locks on *rwlock* (that is, successfully call `pthread_rwlock_rdlock()` *n* times) If so, the thread must perform matching unlocks (that is, it must call `pthread_rwlock_unlock()` *n* times).

Results are undefined if `pthread_rwlock_rdlock()` is called with an uninitialized read-write lock.

If a signal is delivered to a thread waiting for a read-write lock for reading, on return from the signal handler the thread resumes waiting for the read-write lock for reading as if it was not interrupted.

Return Value

If successful, `pthread_rwlock_rdlock()` returns zero. Otherwise, an error number is returned to indicate the error.

EINVAL

The value specified by *attr* or *rwlock* is invalid.

Read Lock With a Nonblocking Read-Write Lock

`pthread_rwlock_tryrdlock(3THR)`

```
#include <pthread.h>
```

```
int pthread_rwlock_tryrdlock(pthread_rwlock_t *rwlock);
```

`pthread_rwlock_tryrdlock(3THR)` applies a read lock as in `pthread_rwlock_rdlock()` with the exception that the function fails if any thread holds a write lock on *rwlock* or there are writers blocked on *rwlock*. (For Solaris threads, see “`rw_tryrdlock(3THR)`” on page 184.)

Return Value

`pthread_rwlock_tryrdlock()` returns zero if the lock for reading on the read-write lock object referenced by *rwlock* is acquired. Otherwise an error number is returned to indicate the error.

EBUSY

The read-write lock could not be acquired for reading because a writer holds the lock or was blocked on it.

Write Lock on Read-Write Lock

`pthread_rwlock_wrlock(3THR)`

```
#include <pthread.h>
```

```
int pthread_rwlock_wrlock(pthread_rwlock_t *rwlock);
```

`pthread_rwlock_wrlock(3THR)` applies a write lock to the read-write lock referenced by *rwlock*. The calling thread acquires the write lock if no other thread (reader or writer) holds the read-write lock *rwlock*. Otherwise, the thread blocks (that is, does not return from the `pthread_rwlock_wrlock()` call) until it can acquire the lock. Results are undefined if the calling thread holds the read-write lock (whether a read or write lock) at the time the call is made.

Implementations are allowed to favor writers over readers to avoid writer starvation. (For instance, the Solaris threads implementation favors writers over readers. See “`rw_wrlock(3THR)`” on page 185.)

Results are undefined if `pthread_rwlock_wrlock()` is called with an uninitialized read-write lock.

If a signal is delivered to a thread waiting for a read-write lock for writing, upon return from the signal handler the thread resumes waiting for the read-write lock for writing as if it was not interrupted.

Return Value

`pthread_rwlock_rwlock()` returns zero if the lock for writing on the read-write lock object referenced by *rwlock* is acquired. Otherwise an error number is returned to indicate the error.

Write Lock With a Nonblocking Read-Write Lock

`pthread_rwlock_trywrlock(3THR)`

```
#include <pthread.h>
```

```
int pthread_rwlock_trywrlock(pthread_rwlock_t *rwlock);
```

`pthread_rwlock_trywrlock(3THR)` applies a write lock like `pthread_rwlock_wrlock()`, with the exception that the function fails if any thread currently holds *rwlock* (for reading or writing). (For Solaris threads, see “`rw_trywrlock(3THR)`” on page 185.)

Results are undefined if `pthread_rwlock_trywrlock()` is called with an uninitialized read-write lock.

If a signal is delivered to a thread waiting for a read-write lock for writing, on return from the signal handler the thread resumes waiting for the read-write lock for writing as if it was not interrupted.

Return Value

If successful, `pthread_rwlock_trywrlock()` returns zero if the lock for writing on the read-write lock object referenced by *rwlock* is acquired. Otherwise, an error number is returned to indicate the error.

EBUSY

The read-write lock could not be acquired for writing because it is already locked for reading or writing.

Unlock a Read-Write Lock

`pthread_rwlock_unlock(3THR)`

```
#include <pthread.h>
```

`pthread_rwlock_unlock(3THR)` releases a lock held on the read-write lock object referenced by *rwlock*. Results are undefined if the read-write lock *rwlock* is not held by the calling thread. (For Solaris threads, see “`rw_unlock(3THR)`” on page 186.)

If `pthread_rwlock_unlock()` is called to release a read lock from the read-write lock object and there are other read locks currently held on this read-write lock object, the read-write lock object remains in the read locked state. If `pthread_rwlock_unlock()` releases the calling thread’s last read lock on this read-write lock object, then the calling thread is no longer one of the owners of the object. If `pthread_rwlock_unlock()` releases the last read lock for this read-write lock object, the read-write lock object will be put in the unlocked state with no owners.

If `pthread_rwlock_unlock()` is called to release a write lock for this read-write lock object, the read-write lock object will be put in the unlocked state with no owners.

If the call to the `pthread_rwlock_unlock()` results in the read-write lock object becoming unlocked and there are multiple threads waiting to acquire the read-write lock object for writing, the scheduling policy is used to determine which thread acquires the read-write lock object for writing. If there are multiple threads waiting to acquire the read-write lock object for reading, the scheduling policy is used to determine the order in which the waiting threads acquire the read-write lock object for reading. If there are multiple threads blocked on *rwlock* for both read locks and write locks, it is unspecified whether the readers acquire the lock first or whether a writer acquires the lock first.

Results are undefined if `pthread_rwlock_unlock()` is called with an uninitialized read-write lock.

Return Value

If successful, `pthread_rwlock_unlock()` returns zero. Otherwise, an error number is returned to indicate the error.

Destroy a Read-Write Lock

`pthread_rwlock_destroy(3THR)`

```
#include <pthread.h>
```

```
int pthread_rwlock_destroy(pthread_rwlock_t *rwlock);
```

```
pthread_rwlock_t rwlock = PTHREAD_RWLOCK_INITIALIZER;
```

`pthread_rwlock_destroy(3THR)` destroys the read-write lock object referenced by *rwlock* and releases any resources used by the lock. The effect of subsequent use of the lock is undefined until the lock is re-initialized by another call to `pthread_rwlock_init()`. An implementation can cause

`pthread_rwlock_destroy()` to set the object referenced by *rwlock* to an invalid value. Results are undefined if `pthread_rwlock_destroy()` is called when any thread holds *rwlock*. Attempting to destroy an uninitialized read-write lock results in undefined behavior. A destroyed read-write lock object can be re-initialized using `pthread_rwlock_init()`; the results of otherwise referencing the read-write lock object after it has been destroyed are undefined. (For Solaris threads, see “`pthread_rwlock_destroy(3THR)`” on page 186.)

Return Value

If successful, `pthread_rwlock_destroy()` returns zero. Otherwise, an error number is returned to indicate the error.

EINVAL

The value specified by *attr* or *rwlock* is invalid.

Synchronization Across Process Boundaries

Each of the synchronization primitives can be set up to be used across process boundaries. This is done quite simply by ensuring that the synchronization variable is located in a shared memory segment and by calling the appropriate `init()` routine, after the primitive has been initialized with its shared attribute set as `interprocess`.

Producer/Consumer Problem Example

Example 4–17 shows the producer/consumer problem with the producer and consumer in separate processes. The main routine maps zero-filled memory (that it shares with its child process) into its address space.

A child process is created that runs the consumer. The parent runs the producer.

This example also shows the drivers for the producer and consumer. The `producer_driver()` simply reads characters from `stdin` and calls `producer()`. The `consumer_driver()` gets characters by calling `consumer()` and writes them to `stdout`.

The data structure for Example 4–17 is the same as that used for the condition variables example (see Example 4–4). Two semaphores represent the number of full and empty buffers and ensure that producers wait until there are empty buffers and that consumers wait until there are full buffers.

EXAMPLE 4-17 Synchronization Across Process Boundaries

```
main() {
    int zfd;
    buffer_t *buffer;
    pthread_mutexattr_t mattr;
    pthread_condattr_t cvattr_less, cvattr_more;

    zfd = open("/dev/zero", O_RDWR);
    buffer = (buffer_t *)mmap(NULL, sizeof(buffer_t),
        PROT_READ|PROT_WRITE, MAP_SHARED, zfd, 0);
    buffer->occupied = buffer->nextin = buffer->nextout = 0;

    pthread_mutex_attr_init(&mattr);
    pthread_mutexattr_setpshared(&mattr,
        PTHREAD_PROCESS_SHARED);

    pthread_mutex_init(&buffer->lock, &mattr);
    pthread_condattr_init(&cvattr_less);
    pthread_condattr_setpshared(&cvattr_less, PTHREAD_PROCESS_SHARED);
    pthread_cond_init(&buffer->less, &cvattr_less);
    pthread_condattr_init(&cvattr_more);
    pthread_condattr_setpshared(&cvattr_more,
        PTHREAD_PROCESS_SHARED);
    pthread_cond_init(&buffer->more, &cvattr_more);

    if (fork() == 0)
        consumer_driver(buffer);
    else
        producer_driver(buffer);
}

void producer_driver(buffer_t *b) {
    int item;

    while (1) {
        item = getchar();
        if (item == EOF) {
            producer(b, '\0');
            break;
        } else
            producer(b, (char)item);
    }
}

void consumer_driver(buffer_t *b) {
    char item;

    while (1) {
        if ((item = consumer(b)) == '\0')
            break;
        putchar(item);
    }
}
```

Interprocess Locking Without the Threads Library

Although not generally recommended, you can do interprocess locking without using the threads library. If this is something you want to do, see the instructions in “Using LWPs Between Processes” on page 205.

Comparing Primitives

The most basic synchronization primitive in threads is the mutual exclusion lock. So, it is the most efficient mechanism in both memory use and execution time. The basic use of a mutual exclusion lock is to serialize access to a resource.

The next most efficient primitive in threads is the condition variable. The basic use of a condition variable is to block on a change of state; that is it provides a thread wait facility. Remember that a mutex lock must be acquired before blocking on a condition variable and must be unlocked after returning from `pthread_cond_wait()`. The mutex lock must also be held across the change of state that occurs before the corresponding call to `pthread_cond_signal()`.

The semaphore uses more memory than the condition variable. It is easier to use in some circumstances because a semaphore variable functions on state rather than on control. Unlike a lock, a semaphore does not have an owner. Any thread can increment a semaphore that has blocked.

The read-write lock permits concurrent reads and exclusive writes to a protected resource. The read-write lock is a single entity that can be locked in *read* or *write* mode. To modify a resource, a thread must first acquire the exclusive write lock. An exclusive write lock is not permitted until all read locks have been released.

Programming With the Operating Environment

This chapter describes how multithreading interacts with the Solaris operating environment and how the operating environment has changed to support multithreading.

- “Process Creation—`exec(2)` and `exit(2)` Issues” on page 139
- “Timers, Alarms, and Profiling” on page 140
- “Nonlocal Goto—`setjmp(3C)` and `longjmp(3C)`” on page 142
- “Resource Limits” on page 142
- “LWPs and Scheduling Classes” on page 142
- “Extending Traditional Signals” on page 145
- “I/O Issues” on page 154

Process Creation—Forking Issues

The default handling of the `fork()` function in the Solaris operating environment is somewhat different from the way `fork()` is handled in POSIX threads, although the Solaris operating environment does support both mechanisms.

Table 5-1 compares the differences and similarities of Solaris and pthreads `fork()` handling. When the comparable interface is not available either in POSIX threads or in Solaris threads, the ‘—’ character appears in the table column.

TABLE 5-1 Comparing POSIX and Solaris `fork()` Handling

	Solaris Operating Environment Interface	POSIX Threads Interface
Fork-one model	<code>fork1(2)</code>	<code>fork(2)</code>
Fork-all model	<code>fork(2)</code>	—

TABLE 5-1 Comparing POSIX and Solaris `fork()` Handling (Continued)

	Solaris Operating Environment Interface	POSIX Threads Interface
Fork safety	—	<code>pthread_atfork(3THR)</code>

The Fork-One Model

As shown in Table 5-1, the behavior of the pthreads `fork(2)` function is the same as that of the Solaris `fork1(2)` function. Both the pthreads `fork(2)` function and the Solaris `fork1(2)` create a new process, duplicating the complete address space in the child, but duplicating only the calling thread in the child process.

This is useful when the child process immediately calls `exec()`, which is what happens after most calls to `fork()`. In this case, the child process does not need a duplicate of any thread other than the one that called `fork()`.

In the child, do not call any library functions after calling `fork()` and before calling `exec()` because one of the library functions might use a lock that was held in the parent at the time of the `fork()`. The child process might execute only Async-Signal-Safe operations until one of the `exec()` handlers is called.

The Fork-One Safety Problem and Solution

In addition to all of the usual concerns such as locking shared data, a library should be well behaved with respect to forking a child process when only one thread is running (the one that called `fork()`). The problem is that the sole thread in the child process might try to grab a lock that is held by a thread that wasn't duplicated in the child.

This is not a problem most programs are likely to encounter. Most programs call `exec()` in the child right after the return from `fork()`. However, if the program wishes to carry out some actions in the child before the call to `exec()`, or never calls `exec()`, then the child *could* encounter deadlock scenarios.

Each library writer should provide a safe solution, although not providing a fork-safe library is not a large concern because this condition is rare.

For example, assume that T1 is in the middle of printing something (and so is holding a lock for `printf()`), when T2 forks a new process. In the child process, if the sole thread (T2) calls `printf()`, it promptly deadlocks.

The POSIX `fork()` or Solaris `fork1()` duplicates only the thread that calls it. (Calling the Solaris `fork()` duplicates all threads, so this issue does not come up.)

To prevent deadlock, ensure that no such locks are being held at the time of forking. The most obvious way to do this is to have the forking thread acquire all the locks that could possibly be used by the child. Because you cannot do this for locks like those in `printf()` (because `printf()` is owned by `libc`), you must ensure that `printf()` is not being used at `fork()` time.

To manage the locks in your library:

- Identify all the locks used by the library.
- Identify the locking order for the locks used by the library. (If a strict locking order is not used, then lock acquisition must be managed carefully.)
- Arrange to acquire those locks at fork time. In Solaris threads this must be done manually, obtaining the locks just before calling `fork1()`, and releasing them right after.

In the following example, the list of locks used by the library is $\{L1, \dots, Ln\}$, and the locking order for these locks is also $L1 \dots Ln$.

```
mutex_lock(L1);
mutex_lock(L2);
fork1(...);
mutex_unlock(L1);
mutex_unlock(L2);
```

In `pthread`s, you can add a call to `pthread_atfork(f1, f2, f3)` in your library's `.init()` section, where `f1()`, `f2()`, `f3()` are defined as follows:

```
f1() /* This is executed just before the process forks. */
{
    mutex_lock(L1); |
    mutex_lock(...); | -- ordered in lock order
    mutex_lock(Ln); |
} v

f2() /* This is executed in the child after the process forks. */
{
    mutex_unlock(L1);
    mutex_unlock(...);
    mutex_unlock(Ln);
}

f3() /* This is executed in the parent after the process forks. */
{
    mutex_unlock(L1);
    mutex_unlock(...);
    mutex_unlock(Ln);
}
```

Another example of deadlock would be a thread in the parent process—other than the one that called Solaris `fork1(2)`—that has locked a mutex. This mutex is copied into the child process in its locked state, but no thread is copied over to unlock the mutex. So, any thread in the child that tries to lock the mutex waits forever.

Virtual Forks—`vfork(2)`

The standard `fork(2)` function is unsafe in multithreaded programs. `vfork(2)` is like `fork(2)` in that only the calling thread is copied in the child process. As in nonthreaded implementations, `vfork()` does not copy the address space for the child process.

Be careful that the thread in the child process does not change memory before it calls `exec(2)`. Remember that `vfork()` gives the parent address space to the child. The parent gets its address space back after the child calls `exec()` or exits. It is important that the child not change the state of the parent.

For example, it is disastrous to create new threads between the call to `vfork()` and the call to `exec()`.

The Solution—`pthread_atfork(3THR)`

Use `pthread_atfork()` to prevent deadlocks whenever you use the fork-one model.

```
#include <pthread.h>

int pthread_atfork(void (*prepare) (void), void (*parent) (void),
                  void (*child) (void) );
```

The `pthread_atfork()` function declares `fork()` handlers that are called before and after `fork()` in the context of the thread that called `fork()`.

- The *prepare* handler is called before `fork()` starts.
- The *parent* handler is called after `fork()` returns in the parent.
- The *child* handler is called after `fork()` returns in the child.

Any one of these can be set to `NULL`. The order in which successive calls to `pthread_atfork()` are made is significant.

For example, a *prepare* handler could acquire all the mutexes needed, and then the *parent* and *child* handlers could release them. This ensures that all the relevant locks are held by the thread that calls the fork function *before* the process is forked, preventing the deadlock in the child.

Using the fork all model avoids the deadlock problem described in “The Fork-One Safety Problem and Solution” on page 136.

Return Values

`pthread_atfork()` returns a zero when it completes successfully. Any other return value indicates that an error occurred. If the following condition is detected, `pthread_atfork(3THR)` fails and returns the corresponding value.

ENOMEM

Insufficient table space exists to record the fork handler addresses.

The Fork-All Model

The Solaris `fork(2)` function duplicates the address space and all the threads (and LWPs) in the child. This is useful, for example, when the child process never calls `exec(2)` but does use its copy of the parent address space. The fork-all functionality is not available in POSIX threads.

Note that when one thread in a process calls Solaris `fork(2)`, threads that are blocked in an interruptible system call return `EINTR`.

Also, be careful not to create locks that are held by both the parent and child processes. This can happen when locks are allocated in memory that is sharable (that is use `mmap()` with the `MAP_SHARED` flag). Note that this is not a problem if the fork-one model is used.

Choosing the Right Fork

You determine whether `fork()` has a “fork-all” or a “fork-one” semantic in your application by linking with the appropriate library. Linking with `-lthread` gives you the “fork-all” semantic for `fork()`, and linking with `-lpthread` gives the “fork-one” semantic for `fork()` (see Figure 7-1 for an explanation of compiling options).

Cautions for Any Fork

Be careful when using global state after a call to any `fork()` function.

For example, when one thread reads a file serially and another thread in the process successfully calls one of the forks, each process then contains a thread that is reading the file. Because the seek pointer for a file descriptor is shared after a `fork()`, the thread in the parent gets some data while the thread in the child gets the other. This introduces gaps in the sequential read accesses.

Process Creation—`exec(2)` and `exit(2)` Issues

Both the `exec(2)` and `exit(2)` system calls work as they do in single-threaded processes except that they destroy all the threads in the address space. Both calls block until all the execution resources (and so all active threads) are destroyed.

When `exec()` rebuilds the process, it creates a single lightweight process (LWP). The process startup code builds the initial thread. As usual, if the initial thread returns, it calls `exit()` and the process is destroyed.

When all the threads in a process exit, the process exits. A call to any `exec()` function from a process with more than one thread terminates all threads, and loads and executes the new executable image. No destructor functions are called.

Timers, Alarms, and Profiling

The “End of Life” announcements for per-LWP timers (see `timer_create(3RT)`) and per-thread alarms (see `alarm(2)` or `setitimer(2)`) were made in the Solaris 2.5 release. Both features are now replaced with the per-process variants described in this section.

Originally, each LWP had a unique realtime interval timer and alarm that a thread bound to the LWP could use. The timer or alarm delivered one signal to the thread when the timer or alarm expired.

Each LWP also had a virtual time or profile interval timer that a thread bound to the LWP could use. When the interval timer expired, either `SIGVTALRM` or `SIGPROF`, as appropriate, was sent to the LWP that owned the interval timer.

Per-LWP POSIX Timers

In the Solaris 2.3 and 2.4 releases, the `timer_create(3RT)` function returned a timer object with a timer ID meaningful only within the calling LWP and with expiration signals delivered to that LWP. Because of this, the only threads that could use the POSIX timer facility were bound threads.

Even with this restricted use, POSIX timers in the Solaris 2.3 and 2.4 releases for multithreaded applications were unreliable about masking the resulting signals and delivering the associated value from the `sigevent` structure.

Beginning with the Solaris 2.5 release, an application that is compiled defining the macro `_POSIX_PER_PROCESS_TIMERS`, or with a value greater than `199506L` for the symbol `_POSIX_C_SOURCE`, can create per-process timers.

Effective with the Solaris 9 Operating Environment, all timers are per-process except for the virtual time and profile interval timers (see `setitimer(2)` for `ITIMER_VIRTUAL` and `ITIMER_PROF`), which remain per-LWP.

The timer IDs of per-process timers are usable from any LWP, and the expiration signals are generated for the process rather than directed to a specific LWP.

The per-process timers are deleted only by `timer_delete(3RT)` or when the process terminates.

Per-Thread Alarms

In the Solaris Operating Environment 2.3 and 2.4 releases, a call to `alarm(2)` or `setitimer(2)` was meaningful only within the calling LWP. Such timers were deleted automatically when the creating LWP terminated. Because of this, the only threads that could use `alarm()` or `setitimer()` were bound threads.

Even with this restricted use, `alarm()` and `setitimer()` timers in Solaris Operating Environment 2.3 and 2.4 multithreaded applications were unreliable about masking the signals from the bound thread that issued these calls. When such masking was not required, then these two system calls worked reliably from bound threads.

With the Solaris Operating Environment 2.5 release, an application linking with `-lpthread` (POSIX) threads got per-process delivery of `SIGALRM` when calling `alarm()`. The `SIGALRM` generated by `alarm()` is generated for the process rather than directed to a specific LWP. Also, the alarm is reset when the process terminates.

Applications compiled with a release before the Solaris Operating Environment 2.5 release, or not linked with `-lpthread`, will continue to see a per-LWP delivery of signals generated by `alarm()` and `setitimer()`

Effective with the Solaris 9 Operating Environment, calls to `alarm()` or to `setitimer(ITIMER_REAL)` will cause the resulting `SIGALRM` signal to be sent to the process.

Profiling

In Solaris releases prior to 2.6, calling `profil()` in a multithreaded program would impact only the calling LWP; the profile state was not inherited at LWP creation time. To profile a multithreaded program with a global profile buffer, each thread needed to issue a call to `profil()` at threads start-up time, and each thread had to be a bound thread. This was cumbersome and did not easily support dynamically turning profiling on and off. In Solaris 2.6 and later releases, the `profil()` system call for multithreaded processes has global impact—that is, a call to `profil()` impacts all LWPs/threads in a process. This may cause applications that depend on the previous per-LWP semantic to break, but it is expected to improve multithreaded programs that wish to turn profiling on and off dynamically at runtime.

Nonlocal Goto—`setjmp(3C)` and `longjmp(3C)`

The scope of `setjmp()` and `longjmp()` is limited to one thread, which is fine most of the time. However, this does mean that a thread that handles a signal can `longjmp()` only when `setjmp()` is performed in the same thread.

Resource Limits

Resource limits are set on the entire process and are determined by adding the resource use of all threads in the process. When a soft resource limit is exceeded, the offending thread is sent the appropriate signal. The sum of the resources used in the process is available through `getrusage(3C)`.

LWPs and Scheduling Classes

The Solaris kernel has three classes of scheduling. The highest-priority scheduling class is Realtime (RT). The middle-priority scheduling class is `system`. The `system` class cannot be applied to a user process. The lowest-priority scheduling class is timeshare (TS), which is also the default class.

Scheduling class is maintained for each LWP. When a process is created, the initial LWP inherits the scheduling class and priority of the creating LWP in the parent process. As more LWPs are created to run unbound threads, they also inherit this scheduling class and priority.

Threads have the scheduling class and priority of their underlying LWPs. Each LWP in a process can have a unique scheduling class and priority that is visible to the kernel. If a thread is bound, it will always be associated with the same LWP.

Thread priorities regulate contention for synchronization objects. By default LWPs are in the timesharing class. For compute-bound multithreading, thread priorities are not very useful. For multithreaded applications that do a lot of synchronization using the MT libraries, thread priorities become more meaningful.

The scheduling class is set by `prionctl(2)`. How you specify the first two arguments determines whether just the calling LWP or all the LWPs of one or more processes are affected. The third argument of `prionctl()` is the command, which can be one of the following.

- `PC_GETCID`. Get the class ID and class attributes for a specific class.
- `PC_GETCLINFO`. Get the class name and class attributes for a specific class.
- `PC_GETPARMS`. Get the class identifier and the class-specific scheduling parameters of a process, an LWP with a process, or a group of processes.
- `PC_SETPARMS`. Set the class identifier and the class-specific scheduling parameters of a process, an LWP with a process, or a group of processes.

Note that `prionctl()` affects the scheduling of the LWP associated with the calling thread. For unbound threads, the calling thread is not guaranteed to be associated with the affected LWP after the call to `prionctl()` returns.

Timeshare Scheduling

Timeshare scheduling distributes the processing resource fairly among the LWPs in this scheduling class. Other parts of the kernel can monopolize the processor for short intervals without degrading response time as seen by the user.

The `prionctl(2)` call sets the `nice(2)` level of one or more processes. The `prionctl()` call also affects the `nice()` level of all the timesharing class LWPs in the process. The `nice()` level ranges from 0 to +20 normally and from -20 to +20 for processes with superuser privilege. The lower the value, the higher the priority.

The dispatch priority of time shared LWPs is calculated from the instantaneous CPU use rate of the LWP and from its `nice()` level. The `nice()` level indicates the relative priority of the LWPs to the timeshare scheduler.

LWPs with a greater `nice()` value get a smaller, but nonzero, share of the total processing. An LWP that has received a larger amount of processing is given lower priority than one that has received little or no processing.

Realtime Scheduling

The Realtime class (RT) can be applied to a whole process or to one or more LWPs in a process. This requires superuser privilege.

Unlike the `nice(2)` level of the timeshare class, LWPs that are classified Realtime can be assigned priorities either individually or jointly. A `prionctl(2)` call affects the attributes of all the Realtime LWPs in the process.

The scheduler always dispatches the highest-priority Realtime LWP. It preempts a lower-priority LWP when a higher-priority LWP becomes runnable. A preempted LWP is placed at the head of its level queue.

A Realtime LWP retains control of a processor until it is preempted, it suspends, or its Realtime priority is changed. LWPs in the RT class have absolute priority over processes in the TS class.

A new LWP inherits the scheduling class of the parent process or LWP. An RT class LWP inherits the parent's time slice, whether finite or infinite.

An LWP with a finite time slice runs until it terminates, blocks (for example, to wait for an I/O event), is preempted by a higher-priority runnable Realtime process, or the time slice expires.

An LWP with an infinite time slice ceases execution only when it terminates, blocks, or is preempted.

Fair Share Scheduling

The fair share scheduler (FSS) scheduling class allows allocation of CPU time based on shares.

By default, the FSS scheduling class uses the same range of priorities (0 to 59) as the TS and interactive (IA) scheduling classes. All LWPs in a process must run in the same scheduling class. The FSS class schedules individual LWPs, not whole processes. Thus, a mix of processes in the FSS and TS/IA classes could result in unexpected scheduling behavior in both cases.

With the use of processor sets, you can mix TS/IA with FSS in one system as long as all the processes running on each processor set are in either the TS/IA or the FSS scheduling class, so they do not compete for the same CPUs.

Fixed Priority Scheduling

The fixed priority scheduling class (FX) assigns fixed priorities and time quantum that are not adjusted to accommodate resource consumption. Process priority can be changed only by the process itself or another appropriately privileged process. For more information about FX, see the `priocntl(1)` and `dispadm(1M)` manual pages.

Threads in this class share the same range of priorities (0 to 59) as the TS and interactive (IA) scheduling classes. TS is usually the default. FX will usually be used in conjunction with TS.

Extending Traditional Signals

The traditional UNIX signal model is extended to threads in a fairly natural way. The key characteristics are that the signal disposition is process-wide, but the signal mask is per-thread. The process-wide disposition of signals is established using the traditional mechanisms (`signal(3C)`, `sigaction(2)`, and so on).

When a signal handler is marked `SIG_DFL` or `SIG_IGN`, the action on receipt of the signal (exit, core dump, stop, continue, or ignore) is performed on the entire receiving process, affecting all threads in the process. For these signals that don't have handlers, the issue of which thread picks the signal is unimportant, because the action on receipt of the signal is carried out on the whole process. See `signal(5)` for basic information about signals.

Each thread has its own signal mask. This lets a thread block some signals while it uses memory or another state that is also used by a signal handler. All threads in a process share the set of signal handlers set up by `sigaction(2)` and its variants.

A thread in one process cannot send a signal to a specific thread in another process. A signal sent by `kill(2)`, `sigsend(2)` or `sigqueue(3RT)` to a process is handled by any one of the receptive threads in the process.

Signals are divided into two categories: traps and exceptions (synchronously generated signals) and interrupts (asynchronously generated signals).

As in traditional UNIX, if a signal is pending, additional occurrences of that signal normally have no additional effect—a pending signal is represented by a bit, not by a counter. However, signals posted via the `sigqueue(3RT)` interface allow multiple instances of the same signal to be queued to the process.

As is the case with single-threaded processes, when a thread receives a signal while blocked in a system call, the thread might return early, either with the `EINTR` error code, or, in the case of I/O calls, with fewer bytes transferred than requested.

Of particular importance to multithreaded programs is the effect of signals on `pthread_cond_wait(3THR)`. This call usually returns without error (a return value of zero) only in response to a `pthread_cond_signal(3THR)` or a `pthread_cond_broadcast(3THR)`, but, if the waiting thread receives a traditional UNIX signal, it returns with a return value of zero even though the wakeup was spurious. The Solaris threads `cond_wait(3THR)` function returns `EINTR` in this circumstance. See “Interrupted Waits on Condition Variables” on page 153 for more information.

Synchronous Signals

Traps (such as `SIGILL`, `SIGFPE`, `SIGSEGV`) result from something a thread does to itself, such as dividing by zero or making reference to nonexistent memory. A trap is handled only by the thread that caused it. Several threads in a process can generate and handle the same type of trap simultaneously.

Extending the idea of signals to individual threads is easy for synchronously-generated signals—the handler is invoked on the thread that generated the synchronous signal.

However, if the process has not chosen to deal with such problems by establishing an appropriate signal handler, then the default action will be taken when a trap occurs, even if the offending thread has the generated signal blocked. The default action for such signals is to terminate the process, perhaps with a core dump.

Because such a synchronous signal usually means that something is seriously wrong with the whole process, and not just with a thread, terminating the process is often a good choice.

Asynchronous Signals

Interrupts (such as `SIGINT` and `SIGIO`) are asynchronous with any thread and result from some action outside the process. They might be signals sent explicitly by another process, or they might represent external actions such as a user typing `Control-c`.

An interrupt can be handled by any thread whose signal mask allows it. When more than one thread is able to receive the interrupt, only one is chosen.

When multiple occurrences of the same signal are sent to a process, then each occurrence can be handled by a separate thread, as long as threads are available that do not have it masked. When all threads have the signal masked, then the signal is marked *pending* and the first thread to unmask the signal handles it.

Continuation Semantics

Continuation semantics are the traditional way to deal with signals. The idea is that when a signal handler returns, control resumes where it was at the time of the interruption. This is well suited for asynchronous signals in single-threaded processes, as shown in Example 5-1.

This is also used as the exception-handling mechanism in some programming languages, such as PL/1.

EXAMPLE 5-1 Continuation Semantics

```
unsigned int nestcount;
```

EXAMPLE 5-1 Continuation Semantics (Continued)

```
unsigned int A(int i, int j) {
    nestcount++;

    if (i==0)
        return(j+1)
    else if (j==0)
        return(A(i-1, 1));
    else
        return(A(i-1, A(i, j-1)));
}

void sig(int i) {
    printf("nestcount = %d\n", nestcount);
}

main() {
    sigset(SIGINT, sig);
    A(4,4);
}
```

Operations on Signals

pthread_sigmask(3THR)

`pthread_sigmask(3THR)` does for a thread what `sigprocmask(2)` does for a process—it sets the thread’s signal mask. When a new thread is created, its initial mask is inherited from its creator.

The call to `sigprocmask()` in a multithreaded process is equivalent to a call to `pthread_sigmask()`. See the `sigprocmask(2)` page for more information.

pthread_kill(3THR)

`pthread_kill(3THR)` is the thread analog of `kill(2)`—it sends a signal to a specific thread. This, of course, is different from sending a signal to a process. When a signal is sent to a process, the signal can be handled by any thread in the process. A signal sent by `pthread_kill()` can be handled only by the specified thread.

Note that you can use `pthread_kill()` to send signals only to threads in the current process. This is because the thread identifier (type `thread_t`) is local in scope—it is not possible to name a thread in any process but your own.

Note also that the action taken (handler, `SIG_DFL`, `SIG_IGN`) on receipt of a signal by the target thread is global, as usual. This means, for example, that if you send `SIGXXX` to a thread, and the `SIGXXX` signal disposition for the process is to kill the process, then the whole process is killed when the target thread receives the signal.

sigwait(2)

For multithreaded programs, `sigwait(2)` is the preferred interface to use, because it deals well with asynchronously generated signals.

`sigwait()` causes the calling thread to wait until any signal identified by its set argument is delivered to the thread. While the thread is waiting, signals identified by the set argument are unmasked, but the original mask is restored when the call returns.

All signals identified by the set argument must be blocked on all threads, including the calling thread; otherwise, `sigwait()` might not work correctly.

Use `sigwait()` to separate threads from asynchronous signals. You can create one thread that is listening for asynchronous signals while your other threads are created to block any asynchronous signals that might be set to this process.

New sigwait() Implementations

Two versions of `sigwait()` are available beginning with the Solaris Operating Environment 2.5 release: the new Solaris Operating Environment 2.5 version, and the POSIX.1c version. New applications and libraries should use the POSIX standard interface, as the Solaris Operating Environment version might not be available in future releases.

The syntax for the two versions of `sigwait()` is shown below.

```
#include <signal.h>

/* the Solaris 2.5 version*/
int sigwait(sigset_t *set);

/* the POSIX.1c version */
int sigwait(const sigset_t *set, int *sig);
```

When the signal is delivered, the POSIX.1c `sigwait()` clears the pending signal and places the signal number in `sig`. Many threads can call `sigwait()` at the same time, but only one thread returns for each signal that is received.

With `sigwait()` you can treat asynchronous signals synchronously—a thread that deals with such signals simply calls `sigwait()` and returns as soon as a signal arrives. By ensuring that all threads (including the caller of `sigwait()`) have such signals masked, you can be sure that signals are handled only by the intended handler and that they are handled safely.

By always masking all signals in all threads, and just calling `sigwait()` as necessary, your application will be much safer for threads that depend on signals.

Usually, you create one or more threads that call `sigwait()` to wait for signals. Because `sigwait()` can retrieve even masked signals, be sure to block the signals of interest in all other threads so they are not accidentally delivered.

When a signal arrives, a thread returns from `sigwait()`, handles the signal, and calls `sigwait()` again to wait for more signals. The signal-handling thread is not restricted to using Async-Signal-Safe functions and can synchronize with other threads in the usual way. (The Async-Signal-Safe category is defined in “MT Interface Safety Levels” on page 160.)

Note – `sigwait()` cannot receive synchronously-generated signals.

sigtimedwait(3RT)

`sigtimedwait(3RT)` is similar to `sigwait(2)` except that it fails and returns an error when a signal is not received in the indicated amount of time.

Thread-Directed Signals

The UNIX signal mechanism is extended with the idea of thread-directed signals. These are just like ordinary asynchronous signals, except that they are sent to a particular thread instead of to a process.

Waiting for asynchronous signals in a separate thread can be safer and easier than installing a signal handler and processing the signals there.

A better way to deal with asynchronous signals is to treat them synchronously. By calling `sigwait(2)`, discussed in “`sigwait(2)`” on page 148, a thread can wait until a signal occurs.

EXAMPLE 5-2 Asynchronous Signals and `sigwait(2)`

```
main() {
    sigset_t set;
    void runA(void);
    int sig;

    sigemptyset(&set);
    sigaddset(&set, SIGINT);
    pthread_sigmask(SIG_BLOCK, &set, NULL);
    pthread_create(NULL, 0, runA, NULL, PTHREAD_DETACHED, NULL);

    while (1) {
        sigwait(&set, &sig);
        printf("nestcount = %d\n", nestcount);
        printf("received signal %d\n", sig);
    }
}

void runA() {
    A(4,4);
}
```

EXAMPLE 5-2 Asynchronous Signals and `sigwait(2)` (Continued)

```
    exit(0);  
}
```

This example modifies the code of Example 5-1: the main routine masks the `SIGINT` signal, creates a child thread that calls the function `A` of the previous example, and then issues `sigwait()` to handle the `SIGINT` signal.

Note that the signal is masked in the compute thread because the compute thread inherits its signal mask from the main thread. The main thread is protected from `SIGINT` while, and only while, it is not blocked inside of `sigwait()`.

Also, note that there is never any danger of having system calls interrupted when you use `sigwait()`.

Completion Semantics

Another way to deal with signals is with completion semantics.

Use completion semantics when a signal indicates that something so catastrophic has happened that there is no reason to continue executing the current code block. The signal handler runs instead of the remainder of the block that had the problem. In other words, the signal handler completes the block.

In Example 5-3, the block in question is the body of the `then` part of the `if` statement. The call to `setjmp(3C)` saves the current register state of the program in `jbuf` and returns 0, thereby executing the block.

EXAMPLE 5-3 Completion Semantics

```
sigjmp_buf jbuf;  
void mult_divide(void) {  
    int a, b, c, d;  
    void problem();  
  
    sigset(SIGFPE, problem);  
    while (1) {  
        if (sigsetjmp(&jbuf) == 0) {  
            printf("Three numbers, please:\n");  
            scanf("%d %d %d", &a, &b, &c);  
            d = a*b/c;  
            printf("%d*d/%d = %d\n", a, b, c, d);  
        }  
    }  
}  
  
void problem(int sig) {  
    printf("Couldn't deal with them, try again\n");  
}
```

EXAMPLE 5-3 Completion Semantics (Continued)

```
    siglongjmp(&jbuf, 1);  
}
```

If a `SIGFPE` (a floating-point exception) occurs, the signal handler is invoked.

The signal handler calls `siglongjmp(3C)`, which restores the register state saved in `jbuf`, causing the program to return from `sigsetjmp()` again (among the registers saved are the program counter and the stack pointer).

This time, however, `sigsetjmp(3C)` returns the second argument of `siglongjmp()`, which is 1. Notice that the block is skipped over, only to be executed during the next iteration of the `while` loop.

Note that you can use `sigsetjmp(3C)` and `siglongjmp(3C)` in multithreaded programs, but be careful that a thread never does a `siglongjmp()` using the results of another thread's `sigsetjmp()`.

Also, `sigsetjmp()` and `siglongjmp()` save and restore the signal mask, but `setjmp(3C)` and `longjmp(3C)` do not.

It is best to use `sigsetjmp()` and `siglongjmp()` when you work with signal handlers.

Completion semantics are often used to deal with exceptions. In particular, the Sun Ada™ programming language uses this model.

Note – Remember, `sigwait(2)` should *never* be used with synchronous signals.

Signal Handlers and Async-Signal Safety

A concept similar to thread safety is Async-Signal safety. Async-Signal-Safe operations are guaranteed not to interfere with operations that are being interrupted.

The problem of Async-Signal safety arises when the actions of a signal handler can interfere with the operation that is being interrupted.

For example, suppose a program is in the middle of a call to `printf(3S)` and a signal occurs whose handler itself calls `printf()`. In this case, the output of the two `printf()` statements would be intertwined. To avoid this, the handler should not call `printf()` itself when `printf()` might be interrupted by a signal.

This problem cannot be solved by using synchronization primitives because any attempted synchronization between the signal handler and the operation being synchronized would produce immediate deadlock.

Suppose that `printf()` is to protect itself by using a mutex. Now suppose that a thread that is in a call to `printf()`, and so holds the lock on the mutex, is interrupted by a signal.

If the handler (being called by the thread that is still inside of `printf()`) itself calls `printf()`, the thread that holds the lock on the mutex will attempt to take it again, resulting in an instant deadlock.

To avoid interference between the handler and the operation, either ensure that the situation never arises (perhaps by masking off signals at critical moments) or invoke only Async-Signal-Safe operations from inside signal handlers.

The only routines that POSIX guarantees to be Async-Signal-Safe are listed in Table 5-2. Any signal handler can safely call in to one of these functions.

TABLE 5-2 Async-Signal-Safe Functions

<code>_exit()</code>	<code>fstat()</code>	<code>read()</code>	<code>sysconf()</code>
<code>access()</code>	<code>getegid()</code>	<code>rename()</code>	<code>tcdrain()</code>
<code>alarm()</code>	<code>geteuid()</code>	<code>rmdir()</code>	<code>tcflow()</code>
<code>cfgetispeed()</code>	<code>getgid()</code>	<code>setgid()</code>	<code>tcflush()</code>
<code>cfgetospeed()</code>	<code>getgroups()</code>	<code>setpgid()</code>	<code>tcgetattr()</code>
<code>cfsetispeed()</code>	<code>getpgrp()</code>	<code>setsid()</code>	<code>tcgetpgrp()</code>
<code>cfsetospeed()</code>	<code>getpid()</code>	<code>setuid()</code>	<code>tcsendbreak()</code>
<code>chdir()</code>	<code>getppid()</code>	<code>sigaction()</code>	<code>tcsetattr()</code>
<code>chmod()</code>	<code>getuid()</code>	<code>sigaddset()</code>	<code>tcsetpgrp()</code>
<code>chown()</code>	<code>kill()</code>	<code>sigdelset()</code>	<code>time()</code>
<code>close()</code>	<code>link()</code>	<code>sigemptyset()</code>	<code>times()</code>
<code>creat()</code>	<code>lseek()</code>	<code>sigfillset()</code>	<code>umask()</code>
<code>dup2()</code>	<code>mkdir()</code>	<code>sigismember()</code>	<code>uname()</code>
<code>dup()</code>	<code>mkfifo()</code>	<code>sigpending()</code>	<code>unlink()</code>
<code>execle()</code>	<code>open()</code>	<code>sigprocmask()</code>	<code>utime()</code>
<code>execve()</code>	<code>pathconf()</code>	<code>sigsuspend()</code>	<code>wait()</code>
<code>fcntl()</code>	<code>pause()</code>	<code>sleep()</code>	<code>waitpid()</code>
<code>fork()</code>	<code>pipe()</code>	<code>stat()</code>	<code>write()</code>

Interrupted Waits on Condition Variables

When an unmasked, caught signal is delivered to a thread while the thread is waiting on a condition variable, then when the signal handler returns, the thread returns from the condition wait with a spurious wakeup (one not caused by a condition signal call from another thread). In this case, the Solaris threads interfaces (`cond_wait()` and `cond_timedwait()`) return `EINTR` while the POSIX threads interfaces (`pthread_cond_wait()` and `pthread_cond_timedwait()`) return 0. In all cases, the associated mutex lock is reacquired before returning from the condition wait.

This does not imply that the mutex is locked while the thread is executing the signal handler. The state of the mutex in the signal handler is undefined.

The implementation of `libthread` in releases of Solaris prior to the Solaris 9 release guaranteed that the mutex was held while in the signal handler. Applications that rely on this old behavior will require revision for Solaris 9 and subsequent releases.

Handler cleanup is illustrated by Example 5-4.

EXAMPLE 5-4 Condition Variables and Interrupted Waits

```
int sig_catcher() {
    sigset_t set;
    void hdlr();

    mutex_lock(&mut);

    sigemptyset(&set);
    sigaddset(&set, SIGINT);
    sigsetmask(SIG_UNBLOCK, &set, 0);

    if (cond_wait(&cond, &mut) == EINTR) {
        /* signal occurred and lock is held */
        cleanup();
        mutex_unlock(&mut);
        return(0);
    }
    normal_processing();
    mutex_unlock(&mut);
    return(1);
}

void hdlr() {
    /* state of the lock is undefined */
    ...
}
```

Assume that the `SIGINT` signal is blocked in all threads on entry to `sig_catcher()` and that `hdlr()` has been established (with a call to `sigaction(2)`) as the handler for the `SIGINT` signal. When an unmasked and caught instance of the `SIGINT` signal

is delivered to the thread while it is in `cond_wait()`, the thread calls `hdlr()`, then returns to the `cond_wait()` function where the lock on the mutex is reacquired, if necessary, and then returns `EINTR` from `cond_wait()`.

Note that whether `SA_RESTART` has been specified as a flag to `sigaction()` has no effect here; `cond_wait(3THR)` is not a system call and is not automatically restarted. When a caught signal occurs while a thread is blocked in `cond_wait()`, the call always returns `EINTR`.

I/O Issues

One of the attractions of multithreaded programming is I/O performance. The traditional UNIX API gave you little assistance in this area—you either used the facilities of the file system or bypassed the file system entirely.

This section shows how to use threads to get more flexibility through I/O concurrency and multibuffering. This section also discusses the differences and similarities between the approaches of synchronous I/O (with threads) and asynchronous I/O (with and without threads).

I/O as a Remote Procedure Call

In the traditional UNIX model, I/O appears to be synchronous, as if you were placing a remote procedure call to the I/O device. Once the call returns, then the I/O has completed (or at least it appears to have completed—a write request, for example, might merely result in the transfer of the data to a buffer in the operating environment).

The advantage of this model is that it is easy to understand because, as a programmer you are very familiar with the concept of procedure calls.

An alternative approach not found in traditional UNIX systems is the asynchronous model, in which an I/O request merely starts an operation. The program must somehow discover when the operation completes.

This approach is not as simple as the synchronous model, but it has the advantage of allowing concurrent I/O and processing in traditional, single-threaded UNIX processes.

Tamed Asynchrony

You can get most of the benefits of asynchronous I/O by using synchronous I/O in a multithreaded program. Where, with asynchronous I/O, you would issue a request and check later to determine when it completes, you can instead have a separate thread perform the I/O synchronously. The main thread can then check (perhaps by calling `pthread_join(3THR)`) for the completion of the operation at some later time.

Asynchronous I/O

In most situations there is no need for asynchronous I/O, since its effects can be achieved with the use of threads, with each thread doing synchronous I/O. However, in a few situations, threads cannot achieve what asynchronous I/O can.

The most straightforward example is writing to a tape drive to make the tape drive stream. Streaming prevents the tape drive from stopping while it is being written to and moves the tape forward at high speed while supplying a constant stream of data that is written to tape.

To do this, the tape driver in the kernel must issue a queued write request when the tape driver responds to an interrupt that indicates that the previous tape-write operation has completed.

Threads cannot guarantee that asynchronous writes will be ordered because the order in which threads execute is indeterminate. Specifying the order of a write to a tape, for example, is not possible.

Asynchronous I/O Operations

```
#include <sys/asynch.h>

int aioread(int fildes, char *bufp, int bufs, off_t offset,
            int whence, aio_result_t *resultp);

int aiowrite(int fildes, const char *bufp, int bufs,
             off_t offset, int whence, aio_result_t *resultp);

aio_result_t *aiowait(const struct timeval *timeout);

int aiocancel(aio_result_t *resultp);
```

`aioread(3AIO)` and `aiowrite(3AIO)` are similar in form to `pread(2)` and `pwrite(2)`, except for the addition of the last argument. Calls to `aioread()` and `aiowrite()` result in the initiation (or queueing) of an I/O operation.

The call returns without blocking, and the status of the call is returned in the structure pointed to by `resultp`. This is an item of type `aio_result_t` that contains the following:

```
int aio_return;
int aio_errno;
```

When a call fails immediately, the failure code can be found in `aio_errno`. Otherwise, this field contains `AIO_INPROGRESS`, meaning that the operation has been successfully queued.

You can wait for an outstanding asynchronous I/O operation to complete by calling `aio_wait(3AIO)`. This returns a pointer to the `aio_result_t` structure supplied with the original `aio_read(3AIO)` or `aio_write(3)` call.

This time `aio_result_t` contains whatever `read(2)` or `write(2)` would have returned if one of them had been called instead of the asynchronous version. If the `read()` or `write()` is successful, `aio_return` contains the number of bytes that were read or written; if it was not successful, `aio_return` is -1, and `aio_errno` contains the error code.

`aio_wait()` takes a *timeout* argument, which indicates how long the caller is willing to wait. As usual, a `NULL` pointer here means that the caller is willing to wait indefinitely, and a pointer to a structure containing a zero value means that the caller is unwilling to wait at all.

You might start an asynchronous I/O operation, do some work, then call `aio_wait()` to wait for the request to complete. Or you can use `SIGIO` to be notified, asynchronously, when the operation completes.

Finally, a pending asynchronous I/O operation can be cancelled by calling `aio_cancel()`. This routine is called with the address of the result area as an argument. This result area identifies which operation is being cancelled.

Shared I/O and New I/O System Calls

When multiple threads are performing I/O operations at the same time with the same file descriptor, you might discover that the traditional UNIX I/O interface is not thread safe. The problem occurs with nonsequential I/O. This uses the `lseek(2)` system call to set the file offset, which is then used in the next `read(2)` or `write(2)` call to indicate where in the file the operation should start. When two or more threads are issuing `lseek()` to the same file descriptor, a conflict results.

To avoid this conflict, use the `pread(2)` and `pwrite(2)` system calls.

```
#include <sys/types.h>
#include <unistd.h>

ssize_t pread(int fildes, void *buf, size_t nbyte, off_t offset);

ssize_t pwrite(int fildes, void *buf, size_t nbyte,
               off_t offset);
```

These behave just like `read(2)` and `write(2)` except that they take an additional argument, the file offset. With this argument, you specify the offset without using `lseek(2)`, so multiple threads can use these routines safely for I/O on the same file descriptor.

Alternatives to `getc(3C)` and `putc(3C)`

An additional problem occurs with standard I/O. Programmers are accustomed to routines such as `getc(3C)` and `putc(3C)` being very quick—they are implemented as macros. Because of this, they can be used within the inner loop of a program with no concerns about efficiency.

However, when they are made thread safe they suddenly become more expensive—they now require (at least) two internal subroutine calls, to lock and unlock a mutex.

To get around this problem, alternative versions of these routines are supplied, `getc_unlocked(3S)` and `putc_unlocked(3C)`.

These do not acquire locks on a mutex and so are as quick as the original, nonthread-safe versions of `getc(3C)` and `putc(3C)`.

However, to use them in a thread-safe way, you must explicitly lock and release the mutexes that protect the standard I/O streams, using `flockfile(3C)` and `funlockfile(3C)`. The calls to these latter routines are placed outside the loop, and the calls to `getc_unlocked()` or `putc_unlocked()` are placed inside the loop.

Safe and Unsafe Interfaces

This chapter defines MT-safety levels for functions and libraries.

- “Thread Safety” on page 159
- “MT Interface Safety Levels” on page 160
- “Async-Signal-Safe Functions” on page 162
- “MT Safety Levels for Libraries” on page 163

Thread Safety

Thread safety is the avoidance of data races—situations in which data are set to either correct or incorrect values, depending upon the order in which multiple threads access and modify the data.

When no sharing is intended, give each thread a private copy of the data. When sharing is important, provide explicit synchronization to make certain that the program behaves in a deterministic manner.

A procedure is thread safe when it is logically correct when executed simultaneously by several threads. At a practical level, it is convenient to recognize three levels of safety.

- Unsafe
- Thread safe - Serializable
- Thread safe - MT-Safe

An unsafe procedure can be made thread safe and serializable by surrounding it with statements to lock and unlock a mutex. Example 6–1 shows three simplified implementations of `fputs()`, initially thread unsafe.

Next is a serializable version of this routine with a single mutex protecting the procedure from concurrent execution problems. Actually, this is stronger synchronization than is usually necessary. When two threads are sending output to different files using `fputs()`, one need not wait for the other—the threads need synchronization only when they are sharing an output file.

The last version is MT-safe. It uses one lock for each file, allowing two threads to print to different files at the same time. So, a routine is MT-safe when it is thread safe and its execution does not negatively affect performance.

EXAMPLE 6-1 Degrees of Thread Safety

```
/* not thread-safe */
fputs(const char *s, FILE *stream) {
    char *p;
    for (p=s; *p; p++)
        putchar((int)*p, stream);
}

/* serializable */
fputs(const char *s, FILE *stream) {
    static mutex_t mut;
    char *p;
    mutex_lock(&mut);
    for (p=s; *p; p++)
        putchar((int)*p, stream);

    mutex_unlock(&mut);
}

/* MT-Safe */
mutex_t m[NFILE];
fputs(const char *s, FILE *stream) {
    static mutex_t mut;
    char *p;
    mutex_lock(&m[fileno(stream)]);
    for (p=s; *p; p++)
        putchar((int)*p, stream);
    mutex_unlock(&m[fileno(stream)]);
}
```

MT Interface Safety Levels

The threads man pages, `man(3THR)`, use the safety level categories listed in Table 6-1 to describe how well an interface supports threads (these categories are explained more fully in the `Intro(3)` reference manual page).

TABLE 6–1 Interface Safety Levels

Category	Description
Safe	This code can be called from a multithreaded application
Safe with exceptions	See the NOTES sections of these pages for a description of the exceptions.
Unsafe	This interface is not safe to use with multithreaded applications unless the application arranges for only one thread at a time to execute within the library.
MT-Safe	This interface is fully prepared for multithreaded access in that it is both <i>safe</i> and it supports some concurrency.
MT-Safe with exceptions	See the NOTES sections of these pages in the <code>man (3THR) : Library Routines</code> for a list of the exceptions.
Async-Signal-Safe	This routine can safely be called from a signal handler. A thread that is executing an Async-Signal-Safe routine does not deadlock with itself when it is interrupted by a signal.
Fork1-Safe	This interface releases locks it has held whenever the Solaris <code>fork1(2)</code> or the POSIX <code>fork(2)</code> is called.

See the section 3 manual pages for the safety levels of library routines.

Some functions have purposely not been made safe for the following reasons.

- Making the interface MT-Safe would have negatively affected the performance of single-threaded applications.
- The library has an Unsafe interface. For example, a function might return a pointer to a buffer in the stack. You can use reentrant counterparts for some of these functions. The reentrant function name is the original function name with “_r” appended.



Caution – There is no way to be certain that a function with a name not ending in “_r” is MT-Safe other than by checking its reference manual page. Use of a function identified as not MT-Safe must be protected by a synchronizing device or by restriction to the initial thread.

Reentrant Functions for Unsafe Interfaces

For most functions with Unsafe interfaces, an MT-Safe version of the routine exists. The name of the new MT-Safe routine is always the name of the old Unsafe routine with “_r” appended. The Table 6–2 “_r” routines are supplied in the Solaris environment.

TABLE 6–2 Reentrant Functions

asctime_r(3c)	gethostbyname_r(3n)	getservbyname_r(3n)
ctermid_r(3s)	gethostent_r(3n)	getservbyport_r(3n)
ctime_r(3c)	getlogin_r(3c)	getservent_r(3n)
fgetgrent_r(3c)	getnetbyaddr_r(3n)	getspent_r(3c)
fgetpwent_r(3c)	getnetbyname_r(3n)	getspnam_r(3c)
fgetspent_r(3c)	getnetent_r(3n)	gmtime_r(3c)
gamma_r(3m)	getnetgrent_r(3n)	lgamma_r(3m)
getauclassent_r(3)	getprotobyname_r(3n)	localtime_r(3c)
getauclassnam_r(3)	getprotobynumber_r(3n)	nis_sperror_r(3n)
getauevent_r(3)	getprotoent_r(3n)	rand_r(3c)
getauevnam_r(3)	getpwent_r(3c)	readdir_r(3c)
getauevnum_r(3)	getpwnam_r(3c)	strtok_r(3c)
getgrent_r(3c)	getpwuid_r(3c)	tmpnam_r(3s)
getgrgid_r(3c)	getrpcbyname_r(3n)	ttyname_r(3c)
getgrnam_r(3c)	getrpcbynumber_r(3n)	
gethostbyaddr_r(3n)	getrpcent_r(3n)	

Async-Signal-Safe Functions

Functions that can safely be called from signal handlers are *Async-Signal-Safe*. The POSIX standard defines and lists Async-Signal-Safe functions (IEEE Std 1003.1-1990, 3.3.1.3 (3)(f), page 55). In addition to the POSIX Async-Signal-Safe functions, these three functions from the Solaris threads library are also Async-Signal-Safe.

- `sema_post(3THR)`
- `thr_sigsetmask(3THR)`, similar to `pthread_sigmask(3THR)`
- `thr_kill(3THR)`, similar to `pthread_kill(3THR)`

MT Safety Levels for Libraries

All routines that can potentially be called by a thread from a multithreaded program should be MT-Safe.

This means that two or more activations of a routine must be able to *correctly* execute concurrently. So, every library interface that a multithreaded program uses must be MT-Safe.

Not all libraries are now MT-Safe. The commonly used libraries that are MT-Safe are listed in Table 6–3. Additional libraries will eventually be modified to be MT-Safe.

TABLE 6–3 Some MT-Safe Libraries

Library	Comments
lib/libc	Interfaces that are not safe have thread-safe interfaces of the form *_r (often with different semantics)
lib/libdl_stubs	To support static switch compiling
lib/libintl	Internationalization library
lib/libm	Math library compliant with System V Interface Definition, Edition 3, X/Open and ANSI C
lib/libmalloc	Space-efficient memory allocation library; see malloc(3X)
lib/libmapmalloc	Alternative mmap(2)-based memory allocation library; see mapmalloc(3X)
lib/libnsl	The TLI interface, XDR, RPC clients and servers, netdir, netselect and getXXbyYY interfaces are not safe, but have thread-safe interfaces of the form getXXbyYY_r
lib/libresolv	Thread-specific errno support
lib/libsocket	Socket library for making network connections
lib/libw	Wide character and wide string functions for supporting multibyte locales
lib/straddr	Network name-to-address translation library
lib/libX11	X11 Windows library routines
lib/libC	C++ runtime shared objects

Unsafe Libraries

Routines in libraries that are not guaranteed to be MT-Safe can safely be called by multithreaded programs only when such calls are single threaded.

Compiling and Debugging

This chapter describes how to compile and debug multithreaded programs.

- “Compiling a Multithreaded Application” on page 165
- “Debugging a Multithreaded Program” on page 169

Compiling a Multithreaded Application

Many options are available for header files, define flags, and linking.

Preparing for Compilation

The following items are required to compile and link a multithreaded program. Except for the C compiler, all should come with your Solaris operating environment.

- A standard C compiler
- Include files:
 - `<thread.h>` and `<pthread.h>`
 - `<errno.h>`, `<limits.h>`, `<signal.h>`, `<unistd.h>`
- The regular Solaris linker, `ln(1)`
- The Solaris threads library (`libthread`), the POSIX threads library (`libpthread`), and possibly the POSIX realtime library (`librt`) for semaphores
- MT-safe libraries (`libc`, `libm`, `libw`, `libintl`, `libnsl`, `libsocket`, `libmalloc`, `libmapmalloc`, and so on)

Choosing Solaris or POSIX Semantics

Certain functions, including the ones listed in Table 7–1, have different semantics in the POSIX 1003.1c standard than in the Solaris 2.4 Operating Environment release, which was based on an earlier POSIX draft. Function definitions are chosen at compile time. See the man Pages(3): Library Routines for a description of the differences in expected parameters and return values.

TABLE 7–1 Functions With POSIX/Solaris Semantic Differences

sigwait(2)	
ctime_r(3C)	asctime_r(3C)
ftrylockfile(3S) - new	getlogin_r(3C)
getgrnam_r(3C)	getgrgid_r(3C)
getpwnam_r(3C)	getpwuid_r(3C)
readdir_r(3C)	ttyname_r(3C)

The Solaris `fork(2)` function duplicates all threads (*fork-all* behavior), while the POSIX `fork(2)` function duplicates only the calling thread (*fork-one* behavior), as does the Solaris `fork1()` function.

Including `<thread.h>` or `<pthread.h>`

The include file `<thread.h>`, used with the `-lthread` library, compiles code that is upward compatible with earlier releases of the Solaris Operating Environment. This library contains both interfaces—those with Solaris semantics and those with POSIX semantics. To call `thr_setconcurrency(3THR)` with POSIX threads, your program needs to include `<thread.h>`.

The include file `<pthread.h>`, used with the `-lpthread` library, compiles code that is conformant with the multithreading interfaces defined by the POSIX 1003.1c standard. For complete POSIX compliance, the define flag `_POSIX_C_SOURCE` should be set to a (long) value ≥ 199506 :

```
cc [flags] file... -D_POSIX_C_SOURCE=N      (where N 199506L)
```

You can mix Solaris threads and POSIX threads in the same application, by including both `<thread.h>` and `<pthread.h>`, and linking with either the `-lthread` or `-lpthread` library.

In mixed use, Solaris semantics prevail when compiling with `-D_REENTRANT` and linking with `-lthread`, whereas POSIX semantics prevail when compiling with `-D_POSIX_C_SOURCE` and linking with `-lpthread`.

Defining `_REENTRANT` or `_POSIX_C_SOURCE`

For POSIX behavior, compile applications with the `-D_POSIX_C_SOURCE` flag set \geq 199506L. For Solaris behavior, compile multithreaded programs with the `-D_REENTRANT` flag. This applies to every module of an application. .

For mixed applications (for example, Solaris threads with POSIX semantics), compile with the `-D_REENTRANT` and `-D_POSIX_PTHREAD_SEMANTICS` flags.

To compile a single-threaded application, define neither the `-D_REENTRANT` nor the `-D_POSIX_C_SOURCE` flag. When these flags are not present, all the old definitions for `errno`, `stdio`, and so on, remain in effect.

Note – Compile single-threaded applications, not linked with either of the thread libraries (`libthread.so.1` or `libpthread.so.1`), without the `-D_REENTRANT` flag. This eliminates performance degradation incurred when macros, such as `putc(3s)`, are converted into reentrant function calls.

To summarize, POSIX applications that define `-D_POSIX_C_SOURCE` get the POSIX 1003.1c semantics for the routines listed in Table 7-1. Applications that define only `-D_REENTRANT` get the Solaris semantics for these routines. Solaris applications that define `-D_POSIX_PTHREAD_SEMANTICS` get the POSIX semantics for these routines, but can still use the Solaris threads interface.

Applications that define both `-D_POSIX_C_SOURCE` and `-D_REENTRANT` get the POSIX semantics.

Linking With `libthread` or `libpthread`

For POSIX threads behavior, load the `libpthread` library. For Solaris threads behavior, load the `libthread` library. Some POSIX programmers might want to link with `-lthread` to preserve the Solaris distinction between `fork()` and `fork1()`. All that `-lpthread` really does is to make `fork()` behave the same way as the Solaris `fork1()` call.

To use `libthread`, specify `-lthread` before `-lc` on the `ld` command line, or last on the `cc` command line.

To use `libpthread`, specify `-lpthread` before `-lc` on the `ld` command line, or last on the `cc` command line.

Prior to Solaris 9, you should not link a nonthreaded program with `-lthread` or `-lpthread`. Doing so establishes multithreading mechanisms at link time that are initiated at runtime. These slow down a single-threaded application, waste system resources, and produce misleading results when you debug your code.

In Solaris 9 and subsequent releases, linking a nonthreaded program with `-lthread` or `-lpthread` makes no semantic difference to the program. No extra threads or LWPs are created and the main (and only) thread executes as a traditional single-threaded process. The only effect on the program is to make system library locks become real locks (as opposed to dummy function calls) and you pay the price of acquiring uncontended locks.

Figure 7-1 summarizes the compile options.

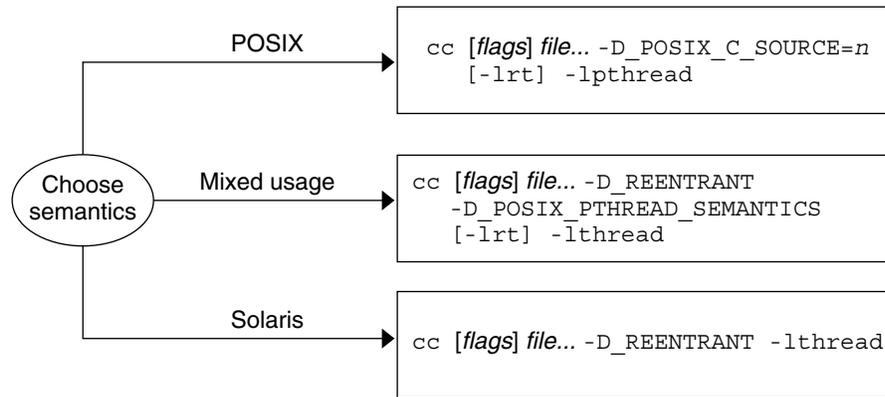


FIGURE 7-1 Compilation Flowchart

In mixed usage, you need to include both `thread.h` and `pthread.h`.

All calls to `libthread` and `libpthread` are no-ops if the application does not link `-lthread` or `-lpthread`. The runtime library `libc` has many predefined `libthread` and `libpthread` stubs that are null procedures. True procedures are interposed by `libthread` or `libpthread` when the application links both `libc` and the thread library.

Note – For C++ programs that use threads, use the `-mt` option, rather than `-lthread`, to compile and link your application. The `-mt` option links with `libthread` and ensures proper library linking order. Using `-lthread` might cause your program to core dump.

Linking With `-lrt` for POSIX Semaphores

The Solaris semaphore routines, `sema_*(3THR)`, are contained in the `libthread` library. By contrast, you link with the `-lrt` library to get the standard `sem_*(3R)` POSIX 1003.1c semaphore routines described in “Semaphores” on page 114.

Link Old With New

Table 7–2 shows that multithreaded object modules should be linked with old object modules only with great caution.

TABLE 7–2 Compiling With and Without the `_REENTRANT` Flag

The File Type	Compiled	Reference	And Return
Old object files (nonthreaded) and new object files	Without the <code>_REENTRANT</code> or <code>_POSIX_C_SOURCE</code> flag	Static storage	The traditional <code>errno</code>
New object files	With the <code>_REENTRANT</code> or <code>_POSIX_C_SOURCE</code> flag	<code>__errno</code> , the new binary entry point	The address of the thread's definition of <code>errno</code>
Programs using TLI in <code>libnsl</code> ¹	With the <code>_REENTRANT</code> or <code>_POSIX_C_SOURCE</code> flag (required)	<code>__t_errno</code> , a new entry point	The address of the thread's definition of <code>t_errno</code> .

1. Include `tuser.h` to get the TLI global error variable.

The Alternate `libthread`

The Solaris 8 Operating Environment introduced an alternate threads library implementation, located in the directories `/usr/lib/lwp` (32-bit) and `/usr/lib/lwp/64` (64-bit). In the Solaris 9 Operating Environment, this implementation is the standard threads implementation found in `/usr/lib` and `/usr/lib/64`.

Debugging a Multithreaded Program

Common Oversights

The following list points out some of the more frequent oversights that can cause bugs in multithreaded programs.

- Passing a pointer to the caller's stack as an argument to a new thread.
- Accessing global memory (shared changeable state) without the protection of a synchronization mechanism.

- Creating deadlocks caused by two threads trying to acquire rights to the same pair of global resources in alternate order (so that one thread controls the first resource and the other controls the second resource and neither can proceed until the other gives up).
- Trying to reacquire a lock already held (recursive deadlock).
- Creating a hidden gap in synchronization protection. This is caused when a code segment protected by a synchronization mechanism contains a call to a function that frees and then reacquires the synchronization mechanism before it returns to the caller. The result is that it appears to the caller that the global data has been protected when it actually has not.
- Mixing UNIX signals with threads—it is better to use the `sigwait(2)` model for handling asynchronous signals.
- Using `setjmp(3C)` and `longjmp(3C)`, and then long-jumping away without releasing the mutex locks.
- Failing to reevaluate the conditions after returning from a call to `*_cond_wait()` or `*_cond_timedwait()`.
- Forgetting that default threads are created `PTHREAD_CREATE_JOINABLE` and must be reclaimed with `pthread_join(3THR)`; note, `pthread_exit(3THR)` does not free up its storage space.
- Making deeply nested, recursive calls and using large automatic arrays can cause problems because multithreaded programs have a more limited stack size than single-threaded programs.
- Specifying an inadequate stack size, or using nondefault stacks.

And, note that multithreaded programs (especially those containing bugs) often behave differently in two successive runs, given identical inputs, because of differences in the thread scheduling order.

In general, multithreading bugs are statistical instead of deterministic. Tracing is usually a more effective method of finding order of execution problems than is breakpoint-based debugging.

Tracing and Debugging With the TNF Utilities

Use the TNF utilities (included as part of the Solaris system) to trace, debug, and gather performance analysis information from your applications and libraries. The TNF utilities integrate trace information from the kernel and from multiple user processes and threads, and so are especially useful for multithreaded code.

With the TNF utilities, you can easily trace and debug multithreaded programs. See the TNF utilities chapter in the *Programming Utilities Guide* for detailed information on using `prex(1)`, `tnfdump(1)`, and other TNF utilities.

Using truss(1)

See `truss(1)` for information on tracing system calls, signals and user-level function calls.

Using mdb(1)

The following `mdb` commands can be used to access the LWPs of a multithreaded program.

TABLE 7-3 MT `mdb` Commands

<code>pid:A</code>	Attaches to process # <i>pid</i> . This stops the process and all its LWPs.
<code>:R</code>	Detaches from process. This resumes the process and all its LWPs.
<code>\$L</code>	Lists all active LWPs in the (stopped) process.
<code>n:l</code>	Switches focus to LWP # <i>n</i> .
<code>\$l</code>	Shows the LWP currently focused.
<code>num:i</code>	Ignores signal number <i>num</i> .

These commands to set conditional breakpoints are often useful.

TABLE 7-4 Setting `mdb` Breakpoints

<code>[label], [count]:b [expression]</code>	Breakpoint is detected when <i>expression</i> equals zero
<code>foo,ffff:b <g7-0xabcdef</code>	Stop at <i>foo</i> when <i>g7</i> = the hex value 0xABCDEF

Using dbx

With the `dbx` utility you can debug and execute source programs written in C++, ANSI C, and FORTRAN. `dbx` accepts the same commands as the Debugger, but uses a standard terminal (TTY) interface. Both `dbx` and the Debugger support debugging multithreaded programs. For a full overview of `dbx` and Debugger features see the `dbx(1)` reference manual page and the *Using Sun Workshop* user's guide.

All the `dbx` options listed in Table 7-5 can support multithreaded applications.

TABLE 7-5 dbx Options for MT Programs

Option	Meaning
cont at line [sig signo id]	Continues execution at <i>line</i> with signal <i>signo</i> . The <i>id</i> , if present, specifies which thread or LWP to continue. The default value is <i>all</i> .
lwp	Displays current LWP. Switches to given LWP [lwpid].
lwps	Lists all LWPs in the current process.
next ... tid	Steps the given thread. When a function call is skipped, all LWPs are implicitly resumed for the duration of that function call. Nonactive threads cannot be stepped.
next ... lid	Steps the given LWP. Does not implicitly resume all LWPs when skipping a function. The LWP on which the given thread is active. Does not implicitly resume all LWP when skipping a function.
step... tid	Steps the given thread. When a function call is skipped, all LWPs are implicitly resumed for the duration of that function call. Nonactive threads cannot be stepped.
step... lid	Steps the given LWP. Does not implicitly resume all LWPs when skipping a function.
stepi... lid	The given LWP.
stepi... tid	The LWP on which the given thread is active.
thread	Displays current thread. Switches to thread <i>tid</i> . In all the following variations, an optional <i>tid</i> implies the current thread.
thread -info [tid]	Prints everything known about the given thread.
thread -locks [tid]	Prints all locks held by the given thread.
thread -suspend [tid]	Puts the given thread into suspended state.
thread -continue [tid]	Unsuspects the given thread.
thread -hide [tid]	<i>Hides</i> the given (or current) thread. It will not appear in the generic <code>threads</code> listing.
thread -unhide [tid]	<i>Unhides</i> the given (or current) thread.
allthread-unhide	<i>Unhides</i> all threads.
threads	Prints the list of all known threads.
threads-all	Prints threads that are not usually printed (zombies).
all filterthreads-mode	Controls whether <code>threads</code> prints all threads or filters them by default.

TABLE 7-5 dbx Options for MT Programs (Continued)

Option	Meaning
auto manualthreads-mode	Enables automatic updating of the thread listing.
threads-mode	Echoes the current modes. Any of the previous forms can be followed by a thread or LWP ID to get the traceback for the specified entity.

Programming With Solaris Threads

This chapter compares the Application Programming Interface (API) for Solaris and POSIX threads, and explains the Solaris features that are not found in POSIX threads.

- “Comparing APIs for Solaris Threads and POSIX Threads” on page 175
- “Unique Solaris Threads Functions” on page 180
- “Similar Synchronization Functions—Read-Write Locks” on page 181
- “Similar Solaris Threads Functions” on page 188
- “Similar Synchronization Functions—Mutual Exclusion Locks” on page 196
- “Similar Synchronization Functions—Condition Variables” on page 199
- “Similar Synchronization Functions—Semaphores” on page 202
- “Special Issues for fork() and Solaris Threads” on page 207

Comparing APIs for Solaris Threads and POSIX Threads

The Solaris threads API and the pthreads API are two solutions to the same problem: building parallelism into application software. Although each API is complete in itself, you can safely mix Solaris threads functions and pthread functions in the same program.

The two APIs do not match exactly, however. Solaris threads supports functions that are not found in pthreads, and pthreads includes functions that are not supported in the Solaris interface. For those functions that *do* match, the associated arguments might not, although the information content is effectively the same.

By combining the two APIs, you can use features not found in one to enhance the other. Similarly, you can run applications using Solaris threads, exclusively, with applications using pthreads, exclusively, on the same system.

Major API Differences

Solaris threads and pthreads are very similar in both API action and syntax. The major differences are listed in Table 8-1.

TABLE 8-1 Unique Solaris Threads and pthreads Features

Solaris Threads (libthread)	POSIX Threads (libpthread)
thr_ prefix for threads function names; sema_ prefix for semaphore function names	pthread_ prefix for pthreads function names; sem_ prefix for semaphore function names
Ability to create “daemon” threads	Cancellation semantics
Suspending and continuing a thread	Scheduling policies

Function Comparison Table

The following table compares Solaris threads functions with pthreads functions. Note that even when Solaris threads and pthreads functions appear to be essentially the same, the arguments to the functions can differ.

When a comparable interface is not available either in pthreads or Solaris threads, a hyphen '-' appears in the column. Entries in the pthreads column that are followed by “POSIX 1003.4” or “POSIX.4” are part of the POSIX Realtime standard specification and are not part of pthreads.

TABLE 8-2 Solaris Threads and POSIX pthreads Comparison

Solaris Threads (libthread)	pthreads (libpthread)
thr_create()	pthread_create()
thr_exit()	pthread_exit()
thr_join()	pthread_join()
thr_yield()	sched_yield() POSIX.4
thr_self()	pthread_self()
thr_kill()	pthread_kill()
thr_sigsetmask()	pthread_sigmask()
thr_setprio()	pthread_setschedparam()
thr_getprio()	pthread_getschedparam()
thr_setconcurrency()	pthread_setconcurrency()

TABLE 8-2 Solaris Threads and POSIX pthreads Comparison *(Continued)*

Solaris Threads (libthread)	pthreads (libpthread)
thr_getconcurrency()	pthread_getconcurrency()
thr_suspend()	-
thr_continue()	-
thr_keycreate()	pthread_key_create()
-	pthread_key_delete()
thr_setspecific()	pthread_setspecific()
thr_getspecific()	pthread_getspecific()
-	pthread_once()
-	pthread_equal()
-	pthread_cancel()
-	pthread_testcancel()
-	pthread_cleanup_push()
-	pthread_cleanup_pop()
-	pthread_setcanceltype()
-	pthread_setcancelstate()
mutex_lock()	pthread_mutex_lock()
mutex_unlock()	pthread_mutex_unlock()
mutex_trylock()	pthread_mutex_trylock()
mutex_init()	pthread_mutex_init()
mutex_destroy()	pthread_mutex_destroy()
cond_wait()	pthread_cond_wait()
cond_timedwait()	pthread_cond_timedwait()
cond_reltimedwait()	pthread_cond_reltimedwait_np()
cond_signal()	pthread_cond_signal()
cond_broadcast()	pthread_cond_broadcast()
cond_init()	pthread_cond_init()
cond_destroy()	pthread_cond_destroy()
rwlock_init()	pthread_rwlock_init()

TABLE 8-2 Solaris Threads and POSIX pthreads Comparison *(Continued)*

Solaris Threads (libthread)	pthreads (libpthread)
<code>rwlock_destroy()</code>	<code>pthread_rwlock_destroy()</code>
<code>rw_rdlock()</code>	<code>pthread_rwlock_rdlock()</code>
<code>rw_wrlock()</code>	<code>pthread_rwlock_wrlock()</code>
<code>rw_unlock()</code>	<code>pthread_rwlock_unlock()</code>
<code>rw_tryrdlock()</code>	<code>pthread_rwlock_tryrdlock()</code>
<code>rw_trywrlock()</code>	<code>pthread_rwlock_trywrlock()</code>
-	<code>pthread_rwlockattr_init()</code>
-	<code>pthread_rwlockattr_destroy()</code>
-	<code>pthread_rwlockattr_getpshared()</code>
-	<code>pthread_rwlockattr_setpshared()</code>
<code>sema_init()</code>	<code>sem_init()</code> POSIX 1003.4
<code>sema_destroy()</code>	<code>sem_destroy()</code> POSIX 1003.4
<code>sema_wait()</code>	<code>sem_wait()</code> POSIX 1003.4
<code>sema_post()</code>	<code>sem_post()</code> POSIX 1003.4
<code>sema_trywait()</code>	<code>sem_trywait()</code> POSIX 1003.4
<code>fork1()</code>	<code>fork()</code>
-	<code>pthread_atfork()</code>
<code>fork()</code> (multiple thread copy)	-
-	<code>pthread_mutexattr_init()</code>
-	<code>pthread_mutexattr_destroy()</code>
type argument in <code>mutex_init()</code>	<code>pthread_mutexattr_setpshared()</code>
-	<code>pthread_mutexattr_getpshared()</code>
-	<code>pthread_mutex_attr_settype()</code>
-	<code>pthread_mutex_attr_gettype()</code>
-	<code>pthread_condattr_init()</code>
-	<code>pthread_condattr_destroy()</code>
type argument in <code>cond_init()</code>	<code>pthread_condattr_setpshared()</code>
-	<code>pthread_condattr_getpshared()</code>

TABLE 8-2 Solaris Threads and POSIX pthreads Comparison *(Continued)*

Solaris Threads (libthread)	pthreads (libpthread)
-	pthread_attr_init()
-	pthread_attr_destroy()
THR_BOUND flag in thr_create()	pthread_attr_setscope()
-	pthread_attr_getscope()
-	pthread_attr_setguardsize()
-	pthread_attr_getguardsize()
stack_size argument in thr_create()	pthread_attr_setstacksize()
-	pthread_attr_getstacksize()
stack_addr argument in thr_create()	pthread_attr_setstackaddr()
-	pthread_attr_getstackaddr()
THR_DETACH flag in thr_create()	pthread_attr_setdetachstate()
-	pthread_attr_getdetachstate()
-	pthread_attr_setschedparam()
-	pthread_attr_getschedparam()
-	pthread_attr_setinheritsched()
-	pthread_attr_getinheritsched()
-	pthread_attr_setschedpolicy()
-	pthread_attr_getschedpolicy()

To use the Solaris threads functions described in this chapter, you must link with the Solaris threads library (-lthread).

Where functionality is virtually the same for both Solaris threads and for pthreads, (even though the function names or arguments might differ), only a brief example consisting of the correct include file and the function prototype is presented. Where return values are not given for the Solaris threads functions, see the appropriate pages in man(3): Library Routines for the function return values.

For more information on Solaris related functions, see the related pthreads documentation for the similarly named function.

Where Solaris threads functions offer capabilities that are not available in pthreads, a full description of the functions is provided.

Unique Solaris Threads Functions

- “Suspend Thread Execution” on page 180
- “Continue a Suspended Thread” on page 181

Suspend Thread Execution

`thr_suspend(3THR)`

`thr_suspend(3THR)` immediately suspends the execution of the thread specified by `target_thread`. On successful return from `thr_suspend()`, the suspended thread is no longer executing.

Once a thread is suspended, subsequent calls to `thr_suspend()` have no effect. Signals cannot awaken the suspended thread; they remain pending until the thread resumes execution.

```
#include <thread.h>
```

```
int thr_suspend(thread_t tid);
```

In the following synopsis, `pthread_t tid` as defined in `pthread` is the same as `thread_t tid` in Solaris threads. `tid` values can be used interchangeably either by assignment or through the use of casts.

```
thread_t tid; /* tid from thr_create() */
```

```
/* pthreads equivalent of Solaris tid from thread created */  
/* with pthread_create() */  
pthread_t ptid;
```

```
int ret;
```

```
ret = thr_suspend(tid);
```

```
/* using pthreads ID variable with a cast */  
ret = thr_suspend((thread_t) ptid);
```

Return Values

`thr_suspend()` returns zero after completing successfully. Any other return value indicates that an error occurred. When the following condition occurs, `thr_suspend()` fails and returns the corresponding value.

ESRCH

tid cannot be found in the current process.

Continue a Suspended Thread

thr_continue(3THR)

thr_continue(3THR) resumes the execution of a suspended thread. Once a suspended thread is continued, subsequent calls to thr_continue() have no effect.

```
#include <thread.h>
```

```
int thr_continue(thread_t tid);
```

A suspended thread will not be awakened by a signal. The signal stays pending until the execution of the thread is resumed by thr_continue().

pthread_t *tid* as defined in pthreads is the same as thread_t *tid* in Solaris threads. *tid* values can be used interchangeably either by assignment or through the use of casts.

```
thread_t tid; /* tid from thr_create() */
```

```
/* pthreads equivalent of Solaris tid from thread created */  
/* with pthread_create() */  
pthread_t ptid;
```

```
int ret;
```

```
ret = thr_continue(tid);
```

```
/* using pthreads ID variable with a cast */  
ret = thr_continue((thread_t) ptid)
```

Return Values

thr_continue() returns zero after completing successfully. Any other return value indicates that an error occurred. When the following condition occurs, thr_continue() fails and returns the corresponding value.

ESRCH

tid cannot be found in the current process.

Similar Synchronization Functions—Read-Write Locks

Read-write locks allow simultaneous read access by many threads while restricting write access to only one thread at a time.

- “Initialize a Read-Write Lock” on page 182
- “Acquire a Read Lock” on page 183
- “Try to Acquire a Read Lock” on page 184
- “Acquire a Write Lock” on page 185
- “Try to Acquire a Write Lock” on page 185
- “Unlock a Read-Write Lock” on page 186
- “Destroy Read-Write Lock State” on page 186

When any thread holds the lock for reading, other threads can also acquire the lock for reading but must wait to acquire the lock for writing. If one thread holds the lock for writing, or is waiting to acquire the lock for writing, other threads must wait to acquire the lock for either reading or writing.

Read-write locks are slower than mutexes, but can improve performance when they protect data that are not frequently written but that are read by many concurrent threads.

Use read-write locks to synchronize threads in this process and other processes by allocating them in memory that is writable and shared among the cooperating processes (see `mmap(2)`) and by initializing them for this behavior.

By default, the acquisition order is not defined when multiple threads are waiting for a read-write lock. However, to avoid writer starvation, the Solaris threads package tends to favor writers over readers.

Read-write locks must be initialized before use.

Initialize a Read-Write Lock

`rwlock_init(3THR)`

```
#include <synch.h> (or #include <thread.h>)
```

```
int rwlock_init(rwlock_t *rwlp, int type, void * arg);
```

Use `rwlock_init(3THR)` to initialize the read-write lock pointed to by `rwlp` and to set the lock state to unlocked. `type` can be one of the following (note that `arg` is currently ignored). (For POSIX threads, see “`pthread_rwlock_init(3THR)`” on page 126.)

- `USYNC_PROCESS` The read-write lock can be used to synchronize threads in this process and other processes. `arg` is ignored.
- `USYNC_THREAD` The read-write lock can be used to synchronize threads in this process, only. `arg` is ignored.

Multiple threads must not initialize the same read-write lock simultaneously. Read-write locks can also be initialized by allocation in zeroed memory, in which case a type of `USYNC_THREAD` is assumed. A read-write lock must not be reinitialized while other threads might be using it.

Initializing Read-Write Locks With Intraprocess Scope

```
#include <thread.h>

rwlock_t rwp;
int ret;

/* to be used within this process only */
ret = rwlock_init(&rwp, USYNC_THREAD, 0);
```

Initializing Read-Write Locks With Interprocess Scope

```
#include <thread.h>

rwlock_t rwp;
int ret;

/* to be used among all processes */
ret = rwlock_init(&rwp, USYNC_PROCESS, 0);
```

Return Values

`rwlock_init()` returns zero after completing successfully. Any other return value indicates that an error occurred. When any of the following conditions occurs, the function fails and returns the corresponding value.

`EINVAL`
Invalid argument.

`EFAULT`
rwp or *arg* points to an illegal address.

Acquire a Read Lock

`rw_rdlock(3THR)`

```
#include <synch.h> (or #include <thread.h>)

int rw_rdlock(rwlock_t *rwp);
```

Use `rw_rdlock(3THR)` to acquire a read lock on the read-write lock pointed to by `rwlp`. When the read-write lock is already locked for writing, the calling thread blocks until the write lock is released. Otherwise, the read lock is acquired. (For POSIX threads, see “`pthread_rwlock_rdlock(3THR)`” on page 127.)

Return Values

`rw_rdlock()` returns zero after completing successfully. Any other return value indicates that an error occurred. When any of the following conditions occurs, the function fails and returns the corresponding value.

`EINVAL`
Invalid argument.

`EFAULT`
`rwlp` points to an illegal address.

Try to Acquire a Read Lock

`rw_tryrdlock(3THR)`

```
#include <synch.h> (or #include <thread.h>)

int rw_tryrdlock(rwlock_t *rwlp);
```

Use `rw_tryrdlock(3THR)` to attempt to acquire a read lock on the read-write lock pointed to by `rwlp`. When the read-write lock is already locked for writing, it returns an error. Otherwise, the read lock is acquired. (For POSIX threads, see “`pthread_rwlock_tryrdlock(3THR)`” on page 128.)

Return Values

`rw_tryrdlock()` returns zero after completing successfully. Any other return value indicates that an error occurred. When any of the following conditions occurs, the function fails and returns the corresponding value.

`EINVAL`
Invalid argument.

`EFAULT`
`rwlp` points to an illegal address.

`EBUSY`
The read-write lock pointed to by `rwlp` was already locked.

Acquire a Write Lock

`rw_wrlck(3THR)`

```
#include <synch.h> (or #include <thread.h>)
```

```
int rw_wrlck(rwlock_t *rwl);
```

Use `rw_wrlck(3THR)` to acquire a write lock on the read-write lock pointed to by `rwl`. When the read-write lock is already locked for reading or writing, the calling thread blocks until all the read locks and write locks are released. Only one thread at a time can hold a write lock on a read-write lock. (For POSIX threads, see “`pthread_rwlock_wrlck(3THR)`” on page 128.)

Return Values

`rw_wrlck()` returns zero after completing successfully. Any other return value indicates that an error occurred. When any of the following conditions occurs, the function fails and returns the corresponding value.

`EINVAL`

Invalid argument.

`EFAULT`

`rwl` points to an illegal address.

Try to Acquire a Write Lock

`rw_trywrlck(3THR)`

```
#include <synch.h> (or #include <thread.h>)
```

```
int rw_trywrlck(rwlock_t *rwl);
```

Use `rw_trywrlck(3THR)` to attempt to acquire a write lock on the read-write lock pointed to by `rwl`. When the read-write lock is already locked for reading or writing, it returns an error. (For POSIX threads, see “`pthread_rwlock_trywrlck(3THR)`” on page 129.)

Return Values

`rw_trywrlck()` returns zero after completing successfully. Any other return value indicates that an error occurred. When any of the following conditions occurs, the function fails and returns the corresponding value.

EINVAL

Invalid argument.

EFAULT

rwlp points to an illegal address.

EBUSY

The read-write lock pointed to by *rwlp* was already locked.

Unlock a Read-Write Lock

`rw_unlock(3THR)`

```
#include <synch.h> (or #include <thread.h>)
```

```
int rw_unlock(rwlock_t *rwlp);
```

Use `rw_unlock(3THR)` to unlock a read-write lock pointed to by *rwlp*. The read-write lock must be locked and the calling thread must hold the lock either for reading or writing. When any other threads are waiting for the read-write lock to become available, one of them is unblocked. (For POSIX threads, see “`pthread_rwlock_unlock(3THR)`” on page 129.)

Return Values

`rw_unlock()` returns zero after completing successfully. Any other return value indicates that an error occurred. When any of the following conditions occurs, the function fails and returns the corresponding value.

EINVAL

Invalid argument.

EFAULT

rwlp points to an illegal address.

Destroy Read-Write Lock State

`rwlock_destroy(3THR)`

```
#include <synch.h> (or #include <thread.h>)
```

```
int rwlock_destroy(rwlock_t *rwlp);
```

Use `rwlock_destroy(3THR)` to destroy any state associated with the read-write lock pointed to by `rwlp`. The space for storing the read-write lock is not freed. (For POSIX threads, see “`pthread_rwlock_destroy(3THR)`” on page 130.)

Return Values

`rwlock_destroy()` returns zero after completing successfully. Any other return value indicates that an error occurred. When any of the following conditions occurs, the function fails and returns the corresponding value.

EINVAL

Invalid argument.

EFAULT

`rwlp` points to an illegal address.

Read-Write Lock Example

Example 8–1 uses a bank account to demonstrate read-write locks. While the program could allow multiple threads to have concurrent read-only access to the account balance, only a single writer is allowed. Note that the `get_balance()` function needs the lock to ensure that the addition of the checking and saving balances occurs atomically.

EXAMPLE 8–1 Read-Write Bank Account

```
rwlock_t account_lock;
float checking_balance = 100.0;
float saving_balance = 100.0;
...
rwlock_init(&account_lock, 0, NULL);
...

float
get_balance() {
    float bal;

    rw_rdlock(&account_lock);
    bal = checking_balance + saving_balance;
    rw_unlock(&account_lock);
    return(bal);
}

void
transfer_checking_to_savings(float amount) {
    rw_wrlock(&account_lock);
    checking_balance = checking_balance - amount;
    saving_balance = saving_balance + amount;
    rw_unlock(&account_lock);
}
```

Similar Solaris Threads Functions

TABLE 8-3 Similar Solaris Threads Functions

Operation	Destination Discussion
Create a thread	"thr_create(3THR)" on page 188
Get the minimal stack size	"thr_min_stack(3THR)" on page 191
Get the thread identifier	"thr_self(3THR)" on page 191
Yield thread execution	"thr_yield(3THR)" on page 192
Send a signal to a thread	"thr_kill(3THR)" on page 192
Access the signal mask of the calling thread	"thr_sigsetmask(3THR)" on page 192
Terminate a thread	"thr_exit(3THR)" on page 192
Wait for thread termination	"thr_join(3THR)" on page 193
Create a thread-specific data key	"thr_keycreate(3THR)" on page 194
Set thread-specific data	"thr_setspecific(3THR)" on page 194
Get thread-specific data	"thr_getspecific(3THR)" on page 194
Set the thread priority	"thr_setprio(3THR)" on page 195
Get the thread priority	"thr_getprio(3THR)" on page 195

Create a Thread

The `thr_create(3THR)` routine is one of the most elaborate of all the Solaris threads library routines.

`thr_create(3THR)`

Use `thr_create(3THR)` to add a new thread of control to the current process. (For POSIX threads, see "`pthread_create(3THR)`" on page 26.)

Note that the new thread does not inherit pending signals, but it does inherit priority and signal masks.

```
#include <thread.h>
```

```
int thr_create(void *stack_base, size_t stack_size,  
              void *(*start_routine) (void *), void *arg, long flags,
```

```
thread_t *new_thread);
```

```
size_t thr_min_stack(void);
```

stack_base—Contains the address for the stack that the new thread uses. If *stack_base* is NULL then `thr_create()` allocates a stack for the new thread with at least *stack_size* bytes.

stack_size—Contains the size, in number of bytes, for the stack that the new thread uses. If *stack_size* is zero, a default size is used. In most cases, a zero value works best. If *stack_size* is not zero, it must be greater than the value returned by `thr_min_stack()`.

There is no general need to allocate stack space for threads. The threads library allocates 1 megabyte of virtual memory for each thread's stack with no swap space reserved. (The library uses the `-MAP_NORESERVE` option of `mmap(2)` to make the allocations.)

start_routine—Contains the function with which the new thread begins execution. When `start_routine()` returns, the thread exits with the exit status set to the value returned by *start_routine* (see “`thr_exit(3THR)`” on page 192).

arg—Can be anything that is described by `void`, which is typically any 4-byte value. Anything larger must be passed indirectly by having the argument point to it.

Note that you can supply only one argument. To get your procedure to take multiple arguments, encode them as one (such as by putting them in a structure).

flags—Specifies attributes for the created thread. In most cases a zero value works best.

The value in *flags* is constructed from the bitwise inclusive OR of the following:

- `THR_SUSPENDED`—Suspends the new thread and does not execute *start_routine* until the thread is started by `thr_continue()`. Use this to operate on the thread (such as changing its priority) before you run it. The termination of a detached thread is ignored.
- `THR_DETACHED`—Detaches the new thread so that its thread ID and other resources can be reused as soon as the thread terminates. Set this when you do not want to wait for the thread to terminate.

Note – When there is no explicit synchronization to prevent it, an unsuspending, detached thread can die and have its thread ID reassigned to another new thread before its creator returns from `thr_create()`.

- `THR_BOUND`—Permanently binds the new thread to an LWP (the new thread is a *bound thread*).

- `THR_DAEMON`—Marks the new thread as a daemon. A daemon thread is always detached (`THR_DAEMON` implies `THR_DETACHED`). The process exits when all nondaemon threads exit. Daemon threads do not affect the process exit status and are ignored when counting the number of thread exits.

A process can exit either by calling `exit()` or by having every thread in the process that was not created with the `THR_DAEMON` flag call `thr_exit(3THR)`. An application, or a library it calls, can create one or more threads that should be ignored (not counted) in the decision of whether to exit. The `THR_DAEMON` flag identifies threads that are not counted in the process exit criterion.

new_thread—Points to a location (when *new_thread* is not `NULL`) where the ID of the new thread is stored when `thr_create()` is successful. The caller is responsible for supplying the storage this argument points to. The ID is valid only within the calling process.

If you are not interested in this identifier, supply a `NULL` value to *new_thread*.

Return Values

`thr_create()` returns zero when it completes successfully. Any other return value indicates that an error occurred. When any of the following conditions is detected, `thr_create()` fails and returns the corresponding value.

`EAGAIN`

A system limit is exceeded, such as when too many LWPs have been created.

`ENOMEM`

Not enough memory was available to create the new thread.

`EINVAL`

stack_base is not `NULL` and *stack_size* is less than the value returned by `thr_min_stack()`.

Stack Behavior

Stack behavior in Solaris threads is generally the same as that in `pthread`s. For more information about stack setup and operation, see “About Stacks” on page 65.

You can get the absolute minimum on stack size by calling `thr_min_stack()`, which returns the amount of stack space required for a thread that executes a null procedure. Useful threads need more than this, so be very careful when reducing the stack size.

You can specify a custom stack in two ways. The first is to supply a `NULL` for the stack location, thereby asking the runtime library to allocate the space for the stack, but to supply the desired size in the `stacksize` parameter to `thr_create()`.

The other approach is to take overall aspects of stack management and supply a pointer to the stack to `thr_create()`. This means that you are responsible not only for stack allocation but also for stack deallocation—when the thread terminates, you must arrange for the disposal of its stack.

When you allocate your own stack, be sure to append a red zone to its end by calling `mprotect(2)`.

Get the Minimal Stack Size

`thr_min_stack(3THR)`

Use `thr_min_stack(3THR)` to get the minimum stack size for a thread.

```
#include <thread.h>
```

```
size_t thr_min_stack(void);
```

`thr_min_stack()` returns the amount of space needed to execute a null thread (a null thread is a thread that is created to execute a null procedure).

A thread that does more than execute a null procedure should allocate a stack size greater than the size of `thr_min_stack()`.

When a thread is created with a user-supplied stack, the user must reserve enough space to run the thread. In a dynamically linked execution environment, it is difficult to know what the thread minimal stack requirements are.

Most users should not create threads with user-supplied stacks. User-supplied stacks exist only to support applications that require complete control over their execution environments.

Instead, users should let the threads library manage stack allocation. The threads library provides default stacks that should meet the requirements of any created thread.

Get the Thread Identifier

`thr_self(3THR)`

Use `thr_self(3THR)` to get the ID of the calling thread. (For POSIX threads, see “`pthread_self(3THR)`” on page 36.)

```
#include <thread.h>
```

```
thread_t thr_self(void);
```

Yield Thread Execution

thr_yield(3THR)

thr_yield(3THR) causes the current thread to yield its execution in favor of another thread with the same or greater priority; otherwise it has no effect. There is no guarantee that a thread calling thr_yield() will do so.

```
#include <thread.h>

void thr_yield(void);
```

Send a Signal to a Thread

thr_kill(3THR)

thr_kill(3THR) sends a signal to a thread. (For POSIX threads, see “pthread_kill(3THR)” on page 39.)

```
#include <thread.h>
#include <signal.h>
int thr_kill(thread_t target_thread, int sig);
```

Access the Signal Mask of the Calling Thread

thr_sigsetmask(3THR)

Use thr_sigsetmask(3THR) to change or examine the signal mask of the calling thread.

```
#include <thread.h>
#include <signal.h>
int thr_sigsetmask(int how, const sigset_t *set, sigset_t *oset);
```

Terminate a Thread

thr_exit(3THR)

Use thr_exit(3THR) to terminate a thread. (For POSIX threads, see “pthread_exit(3THR)” on page 41.)

```
#include <thread.h>

void thr_exit(void *status);
```

Wait for Thread Termination

thr_join(3THR)

Use `thr_join(3THR)` to wait for a thread to terminate. (For POSIX threads, see “`pthread_join(3THR)`” on page 27.)

```
#include <thread.h>

int thr_join(thread_t tid, thread_t *departedid, void **status);
```

Join specific

```
#include <thread.h>

thread_t tid;
thread_t departedid;
int ret;
void *status;

/* waiting to join thread "tid" with status */
ret = thr_join(tid, &departedid, &status);

/* waiting to join thread "tid" without status */
ret = thr_join(tid, &departedid, NULL);

/* waiting to join thread "tid" without return id and status */
ret = thr_join(tid, NULL, NULL);
```

When the `tid` is `(thread_t)0`, then `thread_join()` waits for any undetached thread in the process to terminate. In other words, when no thread identifier is specified, any undetached thread that exits causes `thread_join()` to return.

Join any

```
#include <thread.h>

thread_t tid;
thread_t departedid;
int ret;
void *status;

/* waiting to join any non-detached thread with status */
```

```
ret = thr_join(0, &departedid, &status);
```

By indicating 0 as the thread id in the Solaris `thr_join()`, a join will take place when any non detached thread in the process exits. The `departedid` will indicate the thread ID of the exiting thread.

Create a Thread-Specific Data Key

Except for the function names and arguments, thread specific data is the same for Solaris as it is for POSIX. The synopses for the Solaris functions are given in this section.

`thr_keycreate(3THR)`

`thr_keycreate(3THR)` allocates a key that is used to identify thread-specific data in a process. (For POSIX threads, see “`pthread_key_create(3THR)`” on page 30.)

```
#include <thread.h>

int thr_keycreate(thread_key_t *keyp,
                 void (*destructor) (void *value));
```

Set Thread-Specific Data

`thr_setspecific(3THR)`

`thr_setspecific(3THR)` binds `value` to the thread-specific data key, `key`, for the calling thread. (For POSIX threads, see “`pthread_setspecific(3THR)`” on page 32.)

```
#include <thread.h>

int thr_setspecific(thread_key_t key, void *value);
```

Get Thread-Specific Data

`thr_getspecific(3THR)`

`thr_getspecific(3THR)` stores the current value bound to `key` for the calling thread into the location pointed to by `valuep`. (For POSIX threads, see “`pthread_getspecific(3THR)`” on page 33.)

```
#include <thread.h>

int thr_getspecific(thread_key_t key, void **valuep);
```

Set the Thread Priority

In Solaris threads, if a thread is to be created with a priority other than that of its parent's, it is created in `SUSPEND` mode. While suspended, the thread's priority is modified using the `thr_setprio(3THR)` function call; then it is continued.

A higher priority thread will receive precedence by libthread over lower priority threads with respect to synchronization object contention.

thr_setprio(3THR)

`thr_setprio(3THR)` changes the priority of the thread, specified by *tid*, within the current process to the priority specified by *newprio*. (For POSIX threads, see "pthread_setschedparam(3THR)" on page 38.)

```
#include <thread.h>

int thr_setprio(thread_t tid, int newprio)
```

By default, threads are scheduled based on fixed priorities that range from zero, the least significant, to 127 the most significant.

```
thread_t tid;
int ret;
int newprio = 20;

/* suspended thread creation */
ret = thr_create(NULL, NULL, func, arg, THR_SUSPENDED, &tid);

/* set the new priority of suspended child thread */
ret = thr_setprio(tid, newprio);

/* suspended child thread starts executing with new priority */
ret = thr_continue(tid);
```

Get the Thread Priority

thr_getprio(3THR)

Use `thr_getprio(3THR)` to get the current priority for the thread. Each thread inherits a priority from its creator. `thr_getprio()` stores the current priority, *tid*, in the location pointed to by *newprio*. (For POSIX threads, see "pthread_getschedparam(3THR)" on page 39.)

```
#include <thread.h>

int thr_getprio(thread_t tid, int *newprio)
```

Similar Synchronization Functions—Mutual Exclusion Locks

- “Initialize a Mutex” on page 196
- “Destroy a Mutex” on page 197
- “Acquire a Mutex” on page 198
- “Release a Mutex” on page 198
- “Try to Acquire a Mutex” on page 198

Initialize a Mutex

mutex_init(3THR)

```
#include <synch.h> (or
#include <thread.h>)

int mutex_init(mutex_t *mp, int type, void *arg);
```

Use `mutex_init(3THR)` to initialize the mutex pointed to by `mp`. The `type` can be one of the following (note that `arg` is currently ignored). (For POSIX threads, see “Initialize a Mutex” on page 87.)

- `USYNC_PROCESS` The mutex can be used to synchronize threads in this and other processes.
- `USYNC_PROCESS_ROBUST` The mutex can be used to *robustly* synchronize threads in this and other processes.
- `USYNC_THREAD` The mutex can be used to synchronize threads in this process only.

When a process dies while holding a `USYNC_PROCESS` lock, subsequent requestors of that lock hang. This is a problem for systems which share locks with client processes because the client processes can be abnormally killed. To avoid the problem of hanging on a lock held by a dead process, use `USYNC_PROCESS_ROBUST` to lock the mutex. `USYNC_PROCESS_ROBUST` adds two capabilities:

- In the case of process death, all owned locks held by that process are unlocked.
- The next requestor for any of the locks owned by the dead process receives the lock, but with an error return indicating that the previous owner died while holding the lock..

Mutexes can also be initialized by allocation in zeroed memory, in which case a *type* of `USYNC_THREAD` is assumed.

Multiple threads must not initialize the same mutex simultaneously. A mutex lock must not be reinitialized while other threads might be using it.

Mutexes With Intraprocess Scope

```
#include <thread.h>

mutex_t mp;
int ret;

/* to be used within this process only */
ret = mutex_init(&mp, USYNC_THREAD, 0);
```

Mutexes With Interprocess Scope

```
#include <thread.h>

mutex_t mp;
int ret;

/* to be used among all processes */
ret = mutex_init(&mp, USYNC_PROCESS, 0);
```

Mutexes With Interprocess Scope-Robust

```
#include <thread.h>

mutex_t mp;
int ret;

/* to be used among all processes */
ret = mutex_init(&mp, USYNC_PROCESS_ROBUST, 0);
```

Destroy a Mutex

`mutex_destroy(3THR)`

```
#include <thread.h>

int mutex_destroy (mutex_t *mp);
```

Use `mutex_destroy(3THR)` to destroy any state associated with the mutex pointed to by *mp*. Note that the space for storing the mutex is not freed. (For POSIX threads, see “`pthread_mutex_destroy(3THR)`” on page 93.)

Acquire a Mutex

mutex_lock(3THR)

```
#include <thread.h>

int mutex_lock(mutex_t *mp);
```

Use `mutex_lock(3THR)` to lock the mutex pointed to by `mp`. When the mutex is already locked, the calling thread blocks until the mutex becomes available (blocked threads wait on a prioritized queue). (For POSIX threads, see “`pthread_mutex_lock(3THR)`” on page 89.)

Release a Mutex

mutex_unlock(3THR)

```
#include <thread.h>

int mutex_unlock(mutex_t *mp);
```

Use `mutex_unlock(3THR)` to unlock the mutex pointed to by `mp`. The mutex must be locked and the calling thread must be the one that last locked the mutex (the owner). (For POSIX threads, see “`pthread_mutex_unlock(3THR)`” on page 91.)

Try to Acquire a Mutex

mutex_trylock(3THR)

```
#include <thread.h>

int mutex_trylock(mutex_t *mp);
```

Use `mutex_trylock(3THR)` to attempt to lock the mutex pointed to by `mp`. This function is a nonblocking version of `mutex_lock()`. (For POSIX threads, see “`pthread_mutex_trylock(3THR)`” on page 92.)

Similar Synchronization Functions—Condition Variables

- “Initialize a Condition Variable” on page 199
- “Destroy a Condition Variable” on page 200
- “Wait for a Condition” on page 200
- “Wait for an Absolute Time” on page 201
- “Wait for a Time Interval” on page 201
- “Unblock One Thread” on page 201
- “Unblock All Threads” on page 202

Initialize a Condition Variable

`cond_init(3THR)`

```
#include <thread.h>
```

```
int cond_init(cond_t *cv, int type, int arg);
```

Use `cond_init(3THR)` to initialize the condition variable pointed to by *cv*. The *type* can be one of the following (note that *arg* is currently ignored). (For POSIX threads, see “`pthread_condattr_init(3THR)`” on page 99.)

- `USYNC_PROCESS` The condition variable can be used to synchronize threads in this and other processes. *arg* is ignored.
- `USYNC_THREAD` The condition variable can be used to synchronize threads in this process only. *arg* is ignored.

Condition variables can also be initialized by allocation in zeroed memory, in which case a type of `USYNC_THREAD` is assumed.

Multiple threads must not initialize the same condition variable simultaneously. A condition variable must not be reinitialized while other threads might be using it.

Condition Variables With Intraprocess Scope

```
#include <thread.h>
```

```
cond_t cv;  
int ret;
```

```
/* to be used within this process only */
ret = cond_init(cv, USYNC_THREAD, 0);
```

Condition Variables With Interprocess Scope

```
#include <thread.h>

cond_t cv;
int ret;

/* to be used among all processes */
ret = cond_init(&cv, USYNC_PROCESS, 0);
```

Destroy a Condition Variable

`cond_destroy(3THR)`

```
#include <thread.h>

int cond_destroy(cond_t *cv);
```

Use `cond_destroy(3THR)` to destroy state associated with the condition variable pointed to by `cv`. The space for storing the condition variable is not freed. (For POSIX threads, see “`pthread_condattr_destroy(3THR)`” on page 100.)

Wait for a Condition

`cond_wait(3THR)`

```
#include <thread.h>

int cond_wait(cond_t *cv, mutex_t *mp);
```

Use `cond_wait(3THR)` to atomically release the mutex pointed to by `mp` and to cause the calling thread to block on the condition variable pointed to by `cv`. The blocked thread can be awakened by `cond_signal()`, `cond_broadcast()`, or when interrupted by delivery of a signal or a `fork()`. (For POSIX threads, see “`pthread_cond_wait(3THR)`” on page 104.)

Wait for an Absolute Time

cond_timedwait(3THR)

```
#include <thread.h>
```

```
int cond_timedwait(cond_t *cv, mutex_t *mp, timestruct_t abstime);
```

Use `cond_timedwait(3THR)` as you would use `cond_wait()`, except that `cond_timedwait()` does not block past the time of day specified by *abstime*. (For POSIX threads, see “`pthread_cond_timedwait(3THR)`” on page 107.)

`cond_timedwait()` always returns with the mutex locked and owned by the calling thread even when returning an error.

The `cond_timedwait()` function blocks until the condition is signaled or until the time of day specified by the last argument has passed. The timeout is specified as the time of day so the condition can be retested efficiently without recomputing the time-out value.

Wait for a Time Interval

cond_reltimedwait(3THR)

```
#include <thread.h>
```

```
int cond_reltimedwait(cond_t *cv, mutex_t *mp,
    timestruct_t reltime);
```

Use `cond_reltimedwait(3THR)` as you would use `cond_timedwait()`, except that `cond_reltimedwait()` takes a relative time interval value in its third argument rather than an absolute time of day value. (For POSIX threads see, `pthread_cond_reltimedwait_np(3THR)`.)

`cond_reltimedwait()` always returns with the mutex locked and owned by the calling thread even when returning an error. The `cond_reltimedwait()` function blocks until the condition is signaled or until the time interval specified by the last argument has elapsed.

Unblock One Thread

cond_signal(3THR)

```
#include <thread.h>
```

```
int cond_signal(cond_t *cv);
```

Use `cond_signal(3THR)` to unblock one thread that is blocked on the condition variable pointed to by `cv`. Call this function under protection of the same mutex used with the condition variable being signaled. Otherwise, the condition could be signaled between its test and `cond_wait()`, causing an infinite wait.

Unblock All Threads

`cond_broadcast(3THR)`

```
#include <thread.h>
```

```
int cond_broadcast(cond_t *cv);
```

Use `cond_broadcast(3THR)` to unblock all threads that are blocked on the condition variable pointed to by `cv`. When no threads are blocked on the condition variable then `cond_broadcast()` has no effect.

Similar Synchronization Functions—Semaphores

Semaphore operations are the same in both the Solaris Operating Environment and the POSIX environment. The function name changed from `sema_` in the Solaris Operating Environment to `sem_` in pthreads.

- “Initialize a Semaphore” on page 202
- “Increment a Semaphore” on page 203
- “Block on a Semaphore Count” on page 204
- “Decrement a Semaphore Count” on page 204
- “Destroy the Semaphore State” on page 204

Initialize a Semaphore

`sema_init(3THR)`

```
#include <thread.h>
```

```
int sema_init(sema_t *sp, unsigned int count, int type,
```

```
void *arg);
```

Use `sema_init(3THR)` to initialize the semaphore variable pointed to by `sp` by `count` amount. `type` can be one of the following (note that `arg` is currently ignored).

`USYNC_PROCESS` The semaphore can be used to synchronize threads in this process and other processes. Only one process should initialize the semaphore. `arg` is ignored.

`USYNC_THREAD` The semaphore can be used to synchronize threads in this process, only. `arg` is ignored.

Multiple threads must not initialize the same semaphore simultaneously. A semaphore must not be reinitialized while other threads might be using it.

Semaphores With Intraprocess Scope

```
#include <thread.h>

sema_t sp;
int ret;
int count;
count = 4;

/* to be used within this process only */
ret = sema_init(&sp, count, USYNC_THREAD, 0);
```

Semaphores With Interprocess Scope

```
#include <thread.h>

sema_t sp;
int ret;
int count;
count = 4;

/* to be used among all the processes */
ret = sema_init (&sp, count, USYNC_PROCESS, 0);
```

Increment a Semaphore

`sema_post(3THR)`

```
#include <thread.h>

int sema_post(sema_t *sp);
```

Use `sema_post(3THR)` to atomically increment the semaphore pointed to by `sp`. When any threads are blocked on the semaphore, one is unblocked.

Block on a Semaphore Count

`sema_wait(3THR)`

```
#include <thread.h>

int sema_wait(sema_t *sp);
```

Use `sema_wait(3THR)` to block the calling thread until the count in the semaphore pointed to by *sp* becomes greater than zero, then atomically decrement it.

Decrement a Semaphore Count

`sema_trywait(3THR)`

```
#include <thread.h>

int sema_trywait(sema_t *sp);
```

Use `sema_trywait(3THR)` to atomically decrement the count in the semaphore pointed to by *sp* when the count is greater than zero. This function is a nonblocking version of `sema_wait()`.

Destroy the Semaphore State

`sem_destroy(3THR)`

```
#include <thread.h>

int sema_destroy(sema_t *sp);
```

Use `sem_destroy(3THR)` to destroy any state associated with the semaphore pointed to by *sp*. The space for storing the semaphore is not freed.

Synchronization Across Process Boundaries

Each of the synchronization primitives can be set up to be used across process boundaries. This is done quite simply by ensuring that the synchronization variable is located in a shared memory segment and by calling the appropriate `init` routine with `type` set to `USYNC_PROCESS`.

If this has been done, then the operations on the synchronization variables work just as they do when `type` is `USYNC_THREAD`.

```
mutex_init(&m, USYNC_PROCESS, 0);
rwlock_init(&rw, USYNC_PROCESS, 0);
cond_init(&cv, USYNC_PROCESS, 0);
sema_init(&s, count, USYNC_PROCESS, 0);
```

Using LWPs Between Processes

Using locks and condition variables between processes does not require using the threads library. The recommended approach is to use the threads library interfaces, but when this is not desirable, then the `_lwp_mutex_*` and `_lwp_cond_*` interfaces can be used as follows:

1. Allocate the locks and condition variables as usual in shared memory (either with `shmop(2)` or `mmap(2)`).
2. Then initialize the newly allocated objects appropriately with the `USYNC_PROCESS` type. Because no interface is available to perform the initialization (`_lwp_mutex_init(2)` and `_lwp_cond_init(2)` do not exist), the objects can be initialized using statically allocated and initialized dummy objects.

For example, to initialize `lockp`:

```
lwp_mutex_t *lwp_lockp;
lwp_mutex_t dummy_shared_mutex = SHARED_MUTEX;
/* SHARED_MUTEX is defined in /usr/include/synch.h */
...
...
lwp_lockp = alloc_shared_lock();
*lwp_lockp = dummy_shared_mutex;
```

Similarly, for condition variables:

```
lwp_cond_t *lwp_condp;
lwp_cond_t dummy_shared_cv = SHARED_CV;
/* SHARED_CV is defined in /usr/include/synch.h */
...
```

```

...
lwp_condp = alloc_shared_cv();
*lwp_condp = dummy_shared_cv;

```

Producer/Consumer Problem Example

Example 8-2 shows the producer/consumer problem with the producer and consumer in separate processes. The main routine maps zero-filled memory (that it shares with its child process) into its address space. Note that `mutex_init()` and `cond_init()` must be called because the type of the synchronization variables is `USYNC_PROCESS`.

A child process is created that runs the consumer. The parent runs the producer.

This example also shows the drivers for the producer and consumer. The `producer_driver()` simply reads characters from `stdin` and calls `producer()`. The `consumer_driver()` gets characters by calling `consumer()` and writes them to `stdout`.

The data structure for Example 8-2 is the same as that used for the solution with condition variables (see “Nested Locking With a Singly Linked List” on page 96).

EXAMPLE 8-2 The Producer/Consumer Problem, Using `USYNC_PROCESS`

```

main() {
    int zfd;
    buffer_t *buffer;

    zfd = open("/dev/zero", O_RDWR);
    buffer = (buffer_t *)mmap(NULL, sizeof(buffer_t),
        PROT_READ|PROT_WRITE, MAP_SHARED, zfd, 0);
    buffer->occupied = buffer->nextin = buffer->nextout = 0;

    mutex_init(&buffer->lock, USYNC_PROCESS, 0);
    cond_init(&buffer->less, USYNC_PROCESS, 0);
    cond_init(&buffer->more, USYNC_PROCESS, 0);
    if (fork() == 0)
        consumer_driver(buffer);
    else
        producer_driver(buffer);
}

void producer_driver(buffer_t *b) {
    int item;

    while (1) {
        item = getchar();
        if (item == EOF) {
            producer(b, '\0');
            break;
        } else
            producer(b, (char)item);
    }
}

```

EXAMPLE 8-2 The Producer/Consumer Problem, Using USYNC_PROCESS (Continued)

```
    }  
}  
  
void consumer_driver(buffer_t *b) {  
    char item;  
  
    while (1) {  
        if ((item = consumer(b)) == '\0')  
            break;  
        putchar(item);  
    }  
}
```

A child process is created to run the consumer; the parent runs the producer.

Special Issues for fork() and Solaris Threads

Solaris threads and POSIX threads define the behavior of `fork()` differently. See “Process Creation—`exec(2)` and `exit(2)` Issues” on page 139 for a thorough discussion of `fork()` issues.

Solaris `libthread` supports both `fork()` and `fork1()`. The `fork()` call has “fork-all” semantics—it duplicates everything in the process, including threads and LWPs, creating a true clone of the parent. The `fork1()` call creates a clone that has only one thread; the process state and address space are duplicated, but only the calling thread is cloned.

POSIX `libpthread` supports only `fork()`, which has the same semantics as `fork1()` in Solaris threads.

Whether `fork()` has “fork-all” semantics or “fork-one” semantics is dependent on which library is used. Linking with `-lthread` assigns “fork-all” semantics to `fork()`, while linking with `-lpthread` assigns “fork-one” semantics to `fork()`.

See “Linking With `libthread` or `libpthread`” on page 167 for more details.

Programming Guidelines

This chapter gives some pointers on programming with threads. Most pointers apply to both Solaris and POSIX threads, but where functionality differs, it is noted. Changing from single-threaded thinking to multithreaded thinking is emphasized in this chapter.

- “Rethinking Global Variables” on page 209
- “Providing for Static Local Variables” on page 210
- “Synchronizing Threads” on page 211
- “Avoiding Deadlock” on page 214
- “Following Some Basic Guidelines” on page 216
- “Creating and Using Threads” on page 217
- “Working With Multiprocessors” on page 219
- “Summary” on page 224

Rethinking Global Variables

Historically, most code has been designed for single-threaded programs. This is especially true for most of the library routines called from C programs. The following implicit assumptions were made for single-threaded code:

- When you write into a global variable and then, a moment later, read from it, what you read is exactly what you just wrote.
- This is also true for nonglobal, static storage.
- You do not need synchronization because there is nothing to synchronize with.

The next few examples discuss some of the problems that arise in multithreaded programs because of these assumptions, and how you can deal with them.

Traditional, single-threaded C and UNIX have a convention for handling errors detected in system calls. System calls can return anything as a functional value (for example, `write()` returns the number of bytes that were transferred). However, the value `-1` is reserved to indicate that something went wrong. So, when a system call returns `-1`, you know that it failed.

EXAMPLE 9-1 Global Variables and *errno*

```
extern int errno;
...
if (write(file_desc, buffer, size) == -1) {
    /* the system call failed */
    fprintf(stderr, "something went wrong, "
               "error code = %d\n", errno);
    exit(1);
}
...
```

Rather than return the actual error code (which could be confused with normal return values), the error code is placed into the global variable `errno`. When the system call fails, you can look in `errno` to find out what went wrong.

Now consider what happens in a multithreaded environment when two threads fail at about the same time, but with different errors. Both expect to find their error codes in `errno`, but one copy of `errno` cannot hold both values. This global variable approach simply does not work for multithreaded programs.

Threads solves this problem through a conceptually new storage class—thread-specific data. This storage is similar to global storage in that it can be accessed from any procedure in which a thread might be running. However, it is private to the thread—when two threads refer to the thread-specific data location of the same name, they are referring to two different areas of storage.

So, when using threads, each reference to `errno` is thread specific because each thread has a private copy of `errno`. This is achieved in this implementation by making `errno` a macro that expands to a function call.

Providing for Static Local Variables

Example 9-2 shows a problem similar to the `errno` problem, but involving static storage instead of global storage. The function `gethostbyname(3NSL)` is called with the computer name as its argument. The return value is a pointer to a structure containing the required information for contacting the computer through network communications.

EXAMPLE 9-2 The `gethostbyname()` Problem

```
struct hostent *gethostbyname(char *name) {
    static struct hostent result;
        /* Lookup name in hosts database */
        /* Put answer in result */
    return(&result);
}
```

Returning a pointer to a local variable is generally not a good idea, although it works in this case because the variable is static. However, when two threads call this variable at once with different computer names, the use of static storage conflicts.

Thread-specific data could be used as a replacement for static storage, as in the `errno` problem, but this involves dynamic allocation of storage and adds to the expense of the call.

A better way to handle this kind of problem is to make the caller of `gethostbyname()` supply the storage for the result of the call. This is done by having the caller supply an additional argument, an output argument, to the routine. This requires a new interface to `gethostbyname()`.

This technique is used in threads to fix many of these problems. In most cases, the name of the new interface is the old name with “_r” appended, as in `gethostbyname_r(3NSL)`.

Synchronizing Threads

The threads in an application must cooperate and synchronize when sharing the data and the resources of the process.

A problem arises when multiple threads call something that manipulates an object. In a single-threaded world, synchronizing access to such objects is not a problem, but as Example 9-3 illustrates, this is a concern with multithreaded code. (Note that the `printf(3S)` function is safe to call for a multithreaded program; this example illustrates what could happen if `printf()` were not safe.)

EXAMPLE 9-3 The `printf()` Problem

```
/* thread 1: */
    printf("go to statement reached");

/* thread 2: */
    printf("hello world");
```

EXAMPLE 9-3 The `printf()` Problem (Continued)

```
printed on display:  
  go to hello
```

Single-Threaded Strategy

One strategy is to have a single, application-wide mutex lock that is acquired whenever any thread in the application is running and is released before it must block. Since only one thread can be accessing shared data at any one time, each thread has a consistent view of memory.

Because this is effectively a single-threaded program, very little is gained by this strategy.

Reentrance

A better approach is to take advantage of the principles of modularity and data encapsulation. A reentrant function is one that behaves correctly if it is called simultaneously by several threads. Writing a reentrant function is a matter of understanding just what *behaves correctly* means for this particular function.

Functions that are callable by several threads must be made reentrant. This might require changes to the function interface or to the implementation.

Functions that access global state, like memory or files, have reentrance problems. These functions need to protect their use of global state with the appropriate synchronization mechanisms provided by threads.

The two basic strategies for making functions in modules reentrant are code locking and data locking.

Code Locking

Code locking is done at the function call level and guarantees that a function executes entirely under the protection of a lock. The assumption is that all access to data is done through functions. Functions that share data should execute under the same lock.

Some parallel programming languages provide a construct called a monitor that implicitly does code locking for functions that are defined within the scope of the monitor. A monitor can also be implemented by a mutex lock.

Functions under the protection of the same mutex lock or within the same monitor are guaranteed to execute atomically with respect to each other.

Data Locking

Data locking guarantees that access to a collection of data is maintained consistently. For data locking, the concept of locking code is still there, but code locking is around references to shared (global) data, only. For a mutual exclusion locking protocol, only one thread can be in the critical section for each collection of data.

Alternatively, in a multiple readers, single writer protocol, several readers can be allowed for each collection of data or one writer. Multiple threads can execute in a single module when they operate on different data collections and do not conflict on a single collection for the multiple readers, single writer protocol. So, data locking typically allows more concurrency than does code locking.

What strategy should you use when using locks (whether implemented with mutexes, condition variables, or semaphores) in a program? Should you try to achieve maximum parallelism by locking only when necessary and unlocking as soon as possible (fine-grained locking)? Or should you hold locks for long periods to minimize the overhead of taking and releasing them (coarse-grained locking)?

The granularity of the lock depends on the amount of data it protects. A very coarse-grained lock might be a single lock to protect all data. Dividing how the data is protected by the appropriate number of locks is very important. Too fine a grain of locking can degrade performance. The overhead associated with acquiring and releasing locks can become significant when there are too many locks.

The common wisdom is to start with a coarse-grained approach, identify bottlenecks, and add finer-grained locking where necessary to alleviate the bottlenecks. This is reasonably sound advice, but use your own judgment about finding the balance between maximizing parallelism and minimizing lock overhead.

Invariants

For both code locking and data locking, *invariants* are important to control locking complexity. An invariant is a condition or relation that is always true.

The definition is modified somewhat for concurrent execution: an invariant is a condition or relation that is true when the associated lock is being set. Once the lock is set, the invariant can be false. However, the code holding the lock must reestablish the invariant before releasing the lock.

An invariant can also be a condition or relation that is true when a lock is being set. Condition variables can be thought of as having an invariant that is the condition.

EXAMPLE 9-4 Testing the Invariant With assert(3X)

```
mutex_lock (&lock);
while ((condition) == FALSE)
    cond_wait (&cv, &lock);
assert ((condition) == TRUE);
```

EXAMPLE 9-4 Testing the Invariant With `assert(3X)` (Continued)

```
.  
. .  
.  
mutex_unlock(&lock);
```

The `assert()` statement is testing the invariant. The `cond_wait()` function does not preserve the invariant, which is why the invariant must be reevaluated when the thread returns.

Another example is a module that manages a doubly linked list of elements. For each item on the list a good invariant is the forward pointer of the previous item on the list that should also point to the same thing as the backward pointer of the forward item.

Assume this module uses code-based locking and therefore is protected by a single global mutex lock. When an item is deleted or added the mutex lock is acquired, the correct manipulation of the pointers is made, and the mutex lock is released. Obviously, at some point in the manipulation of the pointers the invariant is false, but the invariant is reestablished before the mutex lock is released.

Avoiding Deadlock

Deadlock is a permanent blocking of a set of threads that are competing for a set of resources. Just because some thread can make progress does not mean that there is not a deadlock somewhere else.

The most common error causing deadlock is *self deadlock* or *recursive deadlock*: a thread tries to acquire a lock it is already holding. Recursive deadlock is very easy to program by mistake.

For example, if a code monitor has every module function grabbing the mutex lock for the duration of the call, then any call between the functions within the module protected by the mutex lock immediately deadlocks. If a function calls some code outside the module which, through some circuitous path, calls back into any method protected by the same mutex lock, then it will deadlock too.

The solution for this kind of deadlock is to avoid calling functions outside the module when you don't know whether they will call back into the module without reestablishing invariants and dropping all module locks before making the call. Of course, after the call completes and the locks are reacquired, the state must be verified to be sure the intended operation is still valid.

An example of another kind of deadlock is when two threads, thread 1 and thread 2, each acquires a mutex lock, A and B, respectively. Suppose that thread 1 tries to acquire mutex lock B and thread 2 tries to acquire mutex lock A. Thread 1 cannot proceed and it is blocked waiting for mutex lock B. Thread 2 cannot proceed and it is blocked waiting for mutex lock A. Nothing can change, so this is a permanent blocking of the threads, and a deadlock.

This kind of deadlock is avoided by establishing an order in which locks are acquired (a *lock hierarchy*). When all threads always acquire locks in the specified order, this deadlock is avoided.

Adhering to a strict order of lock acquisition is not always optimal. When thread 2 has many assumptions about the state of the module while holding mutex lock B, giving up mutex lock B to acquire mutex lock A and then reacquiring mutex lock B in order would cause it to discard its assumptions and reevaluate the state of the module.

The blocking synchronization primitives usually have variants that attempt to get a lock and fail if they cannot, such as `mutex_trylock()`. This allows threads to violate the lock hierarchy when there is no contention. When there is contention, the held locks must usually be discarded and the locks reacquired in order.

Deadlocks Related to Scheduling

Because there is no guaranteed order in which locks are acquired, a problem in threaded programs is that a particular thread never acquires a lock, even though it seems that it should.

This usually happens when the thread that holds the lock releases it, lets a small amount of time pass, and then reacquires it. Because the lock was released, it might seem that the other thread should acquire the lock. But, because nothing blocks the thread holding the lock, it continues to run from the time it releases the lock until it reacquires the lock, and so no other thread is run.

You can usually solve this type of problem by calling `thr_yield(3THR)` just before the call to reacquire the lock. This allows other threads to run and to acquire the lock.

Because the time-slice requirements of applications are so variable, the threads library does not impose any. Use calls to `thr_yield()` to make threads share time as you require.

Locking Guidelines

Here are some simple guidelines for locking.

- Try not to hold locks across long operations like I/O where performance can be adversely affected.

- Don't hold locks when calling a function that is outside the module and that might reenter the module.
- In general, start with a coarse-grained approach, identify bottlenecks, and add finer-grained locking where necessary to alleviate the bottlenecks. Most locks are held for short amounts of time and contention is rare, so fix only those locks that have measured contention.
- When using multiple locks, avoid deadlocks by making sure that all threads acquire the locks in the same order.

Following Some Basic Guidelines

- Know what you are importing and whether it is safe.
A threaded program cannot arbitrarily enter nonthreaded code.
- Threaded code can safely refer to unsafe code only from the initial thread.
This ensures that the static storage associated with the initial thread is used only by that thread.
- Sun-supplied libraries are assumed to be *unsafe* unless explicitly documented as *safe*.
If a reference manual entry does not state explicitly that an interface is *MT-Safe*, you should assume that the interface is *unsafe*.
- Use compilation flags to manage binary incompatible source changes. (See Chapter 7, *Compiling and Debugging*, for complete instructions.)
 - `-D_REENTRANT` enables multithreading with the `-lthread` library.
 - `-D_POSIX_C_SOURCE` with `-lthread` gives POSIX threads behavior.
 - `-D_POSIX_PTHREADS_SEMANTICS` with `-lthread` gives both Solaris threads and pthreads interfaces with a preference given to the POSIX interfaces when the two interfaces conflict.
- When making a library safe for multithreaded use, do not thread global process operations.
Do not change global operations (or actions with global side effects) to behave in a threaded manner. For example, if file I/O is changed to per-thread operation, threads cannot cooperate in accessing files.
For thread-specific behavior, or *thread cognizant* behavior, use thread facilities. For example, when the termination of `main()` should terminate only the thread that is exiting `main()`, the end of `main()` should be:

```
thr_exit();
/*NOTREACHED*/
```

Creating and Using Threads

The threads packages will cache the threads data structure and stacks so that the repetitive creation of threads can be reasonably inexpensive.

However, creating and destroying threads as they are required is usually more expensive than managing a pool of threads that wait for independent work.

A good example of this is an RPC server that creates a thread for each request and destroys it when the reply is delivered, instead of trying to maintain a pool of threads to service requests.

While thread creation has less overhead compared to that of process creation, it is not efficient when compared to the cost of a few instructions. Create threads for processing that lasts at least a couple of thousand machine instructions.

Lightweight Processes

Figure 9-1 illustrates the relationship between LWPs and the user and kernel levels.

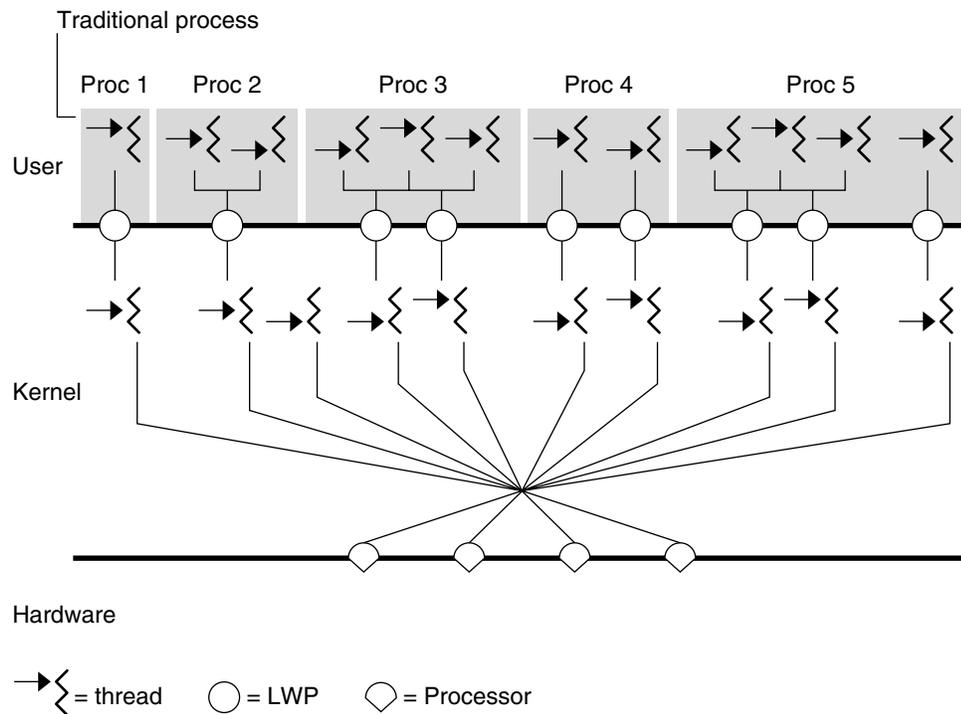


FIGURE 9-1 Multithreading Levels and Relationships

The user-level threads library ensures that the number of LWPs available is adequate for the currently active user-level threads. The operating environment decides which LWP should run on which processor and when. It has no knowledge about user threads. The kernel schedules LWPs onto CPU resources according to their scheduling classes and priorities.

Each LWP is independently dispatched by the kernel, performs independent system calls, incurs independent page faults, and runs in parallel on a multiprocessor system.

An LWP has some capabilities that are not exported directly to threads, such as a special scheduling class.

The new threads library introduced in Solaris 9 actually assigns one LWP to every thread. This is the same as the alternate libthread in Solaris 8.

The new implementation solves many problems that were inherent in the design of the old threads library, principally in the areas of signal handling and concurrency. The new threads library does not have to be told the desired degree of concurrency via `thr_setconcurrency(3THR)` because every thread executes on an LWP.

In future Solaris releases, the threads library might reintroduce multiplexing of unbound threads over LWPs, but with the constraints currently in effect for Solaris 9:

- all runnable threads are attached to LWPs
- no hidden threads are created by the library itself
- a multithreaded process with only one thread has semantics identical to the semantics of a traditional single threaded process.

Unbound Threads

The library invokes LWPs as needed and assigns them to execute runnable threads. The LWP assumes the state of the thread and executes its instructions. If the thread becomes blocked on a synchronization mechanism, the threads library may save the thread state in process memory and assign another thread to the LWP to run.

Bound Threads

Bound threads are guaranteed to execute on the same LWP from the time the thread is created to the time the thread exits.

Thread Creation Guidelines

Here are some simple guidelines for using threads.

- Use threads for independent activities that must do a meaningful amount of work.
- Use bound threads only when a thread needs resources that are available only through the underlying LWP, such as when the thread must be visible to the kernel, as in realtime scheduling.

Working With Multiprocessors

Multithreading lets you take advantage of multiprocessors, primarily through parallelism and scalability. Programmers should be aware of the differences between the memory models of a multiprocessor and a uniprocessor.

Memory consistency is directly interrelated to the processor interrogating memory. For uniprocessors, memory is obviously consistent because there is only one processor viewing memory.

To improve multiprocessor performance, memory consistency is relaxed. You cannot always assume that changes made to memory by one processor are immediately reflected in the other processors' views of that memory.

You can avoid this complexity by using synchronization variables when you use shared or global variables.

Barrier synchronization is sometimes an efficient way to control parallelism on multiprocessors. An example of barriers can be found in Appendix B, Solaris Threads Example: barrier.c.

Another multiprocessor issue is efficient synchronization when threads must wait until all have reached a common point in their execution.

Note – The issues discussed here are not important when the threads synchronization primitives are always used to access shared memory locations.

The Underlying Architecture

When threads synchronize access to shared storage locations using the threads synchronization routines, the effect of running a program on a shared-memory multiprocessor is identical to the effect of running the program on a uniprocessor.

However, in many situations a programmer might be tempted to take advantage of the multiprocessor and use “tricks” to avoid the synchronization routines. As Example 9–5 and Example 9–6 show, such tricks can be dangerous.

Understanding the memory models supported by common multiprocessor architectures helps to understand the dangers.

The major multiprocessor components are:

- The processors themselves
- Store buffers, which connect the processors to their caches
- *Caches*, which hold the contents of recently accessed or modified storage locations
- memory, which is the primary storage (and is shared by all processors).

In the simple traditional model, the multiprocessor behaves as if the processors are connected directly to memory: when one processor stores into a location and another immediately loads from the same location, the second processor loads what was stored by the first.

Caches can be used to speed the average memory access, and the desired semantics can be achieved when the caches are kept consistent with one another.

A problem with this simple approach is that the processor must often be delayed to make certain that the desired semantics are achieved. Many modern multiprocessors use various techniques to prevent such delays, which, unfortunately, change the semantics of the memory model.

Two of these techniques and their effects are explained in the next two examples.

“Shared-Memory” Multiprocessors

Consider the purported solution to the producer/consumer problem shown in Example 9-5.

Although this program works on current SPARC-based multiprocessors, it assumes that all multiprocessors have strongly ordered memory. This program is therefore not portable.

EXAMPLE 9-5 The Producer/Consumer Problem—Shared Memory Multiprocessors

```
char buffer[BFSIZE];
unsigned int in = 0;
unsigned int out = 0;

void producer(char item) {
    do
        /* nothing */
    while
        (in - out == BFSIZE);

    buffer[in%BFSIZE] = item;
    in++;
}

char consumer(void) {
    char item;
    do
        /* nothing */
    while
        (in - out == 0);

    item = buffer[out%BFSIZE];
    out++;
}
```

When this program has exactly one producer and exactly one consumer and is run on a shared-memory multiprocessor, it appears to be correct. The difference between *in* and *out* is the number of items in the buffer.

The producer waits (by repeatedly computing this difference) until there is room for a new item, and the consumer waits until there is an item in the buffer.

For memory that is *strongly ordered* (for instance, a modification to memory on one processor is immediately available to the other processors), this solution is correct (it is correct even taking into account that *in* and *out* will eventually overflow, as long as BFSIZE is less than the largest integer that can be represented in a word).

Shared-memory multiprocessors do not necessarily have strongly ordered memory. A change to memory by one processor is not necessarily available immediately to the other processors. When two changes to different memory locations are made by one processor, the other processors do not necessarily detect the changes in the order in which they were made because changes to memory do not happen immediately.

First the changes are stored in *store buffers* that are not visible to the cache.

The processor checks these store buffers to ensure that a program has a consistent view, but because store buffers are not visible to other processors, a write by one processor does not become visible until it is written to cache.

The synchronization primitives (see Chapter 4, *Programming With Synchronization Objects*) use special instructions that flush the store buffers to cache. So, using locks around your shared data ensures memory consistency.

When memory ordering is very relaxed, Example 9–5 has a problem because the consumer might see that *in* has been incremented by the producer before it sees the change to the corresponding buffer slot.

This is called *weak ordering* because stores made by one processor can appear to happen out of order by another processor (memory, however, is always consistent from the same processor). To fix this, the code should use mutexes to flush the cache.

The trend is toward relaxing memory order. Because of this, programmers are becoming increasingly careful to use locks around all global or shared data.

As demonstrated by Example 9–5 and Example 9–6, locking is essential.

Peterson's Algorithm

The code in Example 9–6 is an implementation of Peterson's Algorithm, which handles mutual exclusion between two threads. This code tries to guarantee that there is never more than one thread in the critical section and that, when a thread calls `mut_excl()`, it enters the critical section sometime "*soon*."

An assumption here is that a thread exits fairly quickly after entering the critical section.

EXAMPLE 9–6 Mutual Exclusion for Two Threads?

```
void mut_excl(int me /* 0 or 1 */) {
    static int loser;
    static int interested[2] = {0, 0};
    int other; /* local variable */

    other = 1 - me;
    interested[me] = 1;
    loser = me;
    while (loser == me && interested[other])
        ;

    /* critical section */
    interested[me] = 0;
}
```

This algorithm works some of the time when it is assumed that the multiprocessor has strongly ordered memory.

Some multiprocessors, including some SPARC-based multiprocessors, have store buffers. When a thread issues a store instruction, the data is put into a store buffer. The buffer contents are eventually sent to the cache, but not necessarily right away. (Note that the caches on each of the processors maintain a consistent view of memory, but modified data does not reach the cache right away.)

When multiple memory locations are stored into, the changes reach the cache (and memory) in the correct order, but possibly after a delay. SPARC-based multiprocessors with this property are said to have total store order (TSO).

When one processor stores into location A and then loads from location B, and another processor stores into location B and loads from location A, the expectation is that either the first processor fetches the newly modified value in location B or the second processor fetches the newly modified value in location A, or both. However, the case in which both processors load the old values simply cannot happen.

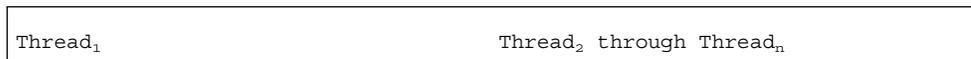
Moreover, with the delays caused by load and store buffers, the “impossible case” can happen.

What could happen with Peterson’s algorithm is that two threads running on separate processors each stores into its own slot of the particular array and then loads from the other slot. They both read the old values (0), assume that the other party is not present, and both enter the critical section. (Note that this is the sort of problem that might not occur when you test a program, but only much later.)

This problem is avoided when you use the threads synchronization primitives, whose implementations issue special instructions to force the writing of the store buffers to the cache.

Parallelizing a Loop on a Shared-Memory Parallel Computer

In many applications, and especially numerical applications, while part of the algorithm can be parallelized, other parts are inherently sequential, as shown in the following:



<pre> while(many_iterations) { sequential_computation --- Barrier --- parallel_computation } </pre>	<pre> while(many_iterations) { --- Barrier --- parallel_computation } </pre>
---	--

For example, you might produce a set of matrixes with a strictly linear computation, then perform operations on the matrixes using a parallel algorithm, then use the results of these operations to produce another set of matrixes, then operate on them in parallel, and so on.

The nature of the parallel algorithms for such a computation is that little synchronization is required during the computation, but synchronization of all the threads employed is required to ensure that the sequential computation is finished before the parallel computation begins.

The barrier forces all the threads that are doing the parallel computation to wait until all threads involved have reached the barrier. When they've reached the barrier, they are released and begin computing together.

Summary

This guide has covered a wide variety of important threads programming issues. Look in Appendix A, Sample Application—Multithreaded grep, for a pthreads program example that uses many of the features and styles that have been discussed. Look in Appendix B, Solaris Threads Example, for a program example that uses Solaris threads.

Further Reading

For more in-depth information about multithreading, see the following book:

- *Programming with Threads* by Steve Kleiman, Devang Shah, and Bart Smaalders (Prentice-Hall, published in 1995)

Sample Application—Multithreaded grep

Description of `tgrep`

The `tgrep` sample program is a multithreaded version of `find(1)` combined with `grep(1)`. `tgrep` supports all but the `-w` (word search) options of the normal `grep`, and a few exclusively available options.

By default, the `tgrep` searches are like the following command:

```
find . -exec grep [options] pattern {} \;
```

For large directory hierarchies, `tgrep` gets results more quickly than the `find` command, depending on the number of processors available. On uniprocessor machines it is about twice as fast, and on four processor machines it is about four times as fast.

The `-e` option changes the way `tgrep` interprets the pattern string. Ordinarily (without the `-e` option) `tgrep` uses a literal string match. With the `-e` option, `tgrep` uses an MT-Safe public domain version of a regular expression handler. The regular expression method is slower.

The `-B` option tells `tgrep` to use the value of the environment variable called `TGLIMIT` to limit the number of threads it will use during a search. This option has no affect if `TGLIMIT` is not set. Because `tgrep` can use a lot of system resources, this is a way to run it politely on a timesharing system.

Getting Online Source Code

Source for `tgrep` is included on the Catalyst Developer's CD. Contact your sales representative to find out how you can get a copy.

Only the multithreaded `main.c` module appears here. Other modules, including those for regular expression handling, plus documentation and Makefiles, are available on the Catalyst Developer's CD.

EXAMPLE A-1 Source Code for `tgrep` Program

```
/* Copyright (c) 1993, 1994 Ron Winacott */
/* This program may be used, copied, modified, and redistributed freely */
/* for ANY purpose, so long as this notice remains intact. */

#define _REENTRANT

#include <stdio.h>
#include <string.h>
#include <stdlib.h>
#include <unistd.h>
#include <assert.h>
#include <errno.h>
#include <ctype.h>
#include <sys/types.h>
#include <time.h>
#include <sys/stat.h>
#include <dirent.h>

#include "version.h"

#include <fcntl.h>
#include <sys/uio.h>
#include <pthread.h>
#include <sched.h>

#ifdef MARK
#include <prof.h> /* to turn on MARK(), use -DMARK to compile (see man prof5)*/
#endif

#include "pmatch.h"

#define PATH_MAX          1024 /* max # of characters in a path name */
#define HOLD_FDS          6 /* stdin,out,err and a buffer */
#define UNLIMITED         99999 /* The default tglimit */
#define MAXREGEXP         10 /* max number of -e options */

#define FB_BLOCK          0x00001
#define FC_COUNT          0x00002
#define FH_HOLDNAME      0x00004
#define FI_IGNCASE        0x00008
```

EXAMPLE A-1 Source Code for tgrep Program (Continued)

```

#define FL_NAMEONLY          0x00010
#define FN_NUMBER            0x00020
#define FS_NOERROR          0x00040
#define FV_REVERSE          0x00080
#define FW_WORD              0x00100
#define FR_RECUR            0x00200
#define FU_UNSORT           0x00400
#define FX_STDIN            0x00800
#define TG_BATCH            0x01000
#define TG_FILEPAT          0x02000
#define FE_REGEXP           0x04000
#define FS_STATS            0x08000
#define FC_LINE             0x10000
#define TG_PROGRESS        0x20000

#define FILET                1
#define DIRT                  2

typedef struct work_st {
    char      *path;
    int       tp;
    struct work_st *next;
} work_t;

typedef struct out_st {
    char      *line;
    int       line_count;
    long      byte_count;
    struct out_st *next;
} out_t;

#define ALPHASIZ             128
typedef struct bm_pattern { /* Boyer - Moore pattern */
    short      p_m; /* length of pattern string */
    short      p_r[ALPHASIZ]; /* "r" vector */
    short      *p_R; /* "R" vector */
    char      *p_pat; /* pattern string */
} BM_PATTERN;

/* bmpmatch.c */
extern BM_PATTERN *bm_makepat(char *p);
extern char *bm_pmatch(BM_PATTERN *pat, register char *s);
extern void bm_freepat(BM_PATTERN *pattern);
BM_PATTERN      *bm_pat; /* the global target read only after main */

/* pmatch.c */
extern char *pmatch(register PATTERN *pattern, register char *string, int *len);
extern PATTERN *makepat(char *string, char *metas);
extern void freepat(register PATTERN *pat);
extern void printpat(PATTERN *pat);
PATTERN      *pm_pat[MAXREGEXP]; /* global targets read only for pmatch */

#include "proto.h" /* function prototypes of main.c */

```

EXAMPLE A-1 Source Code for `tgrep` Program (Continued)

```
/* local functions to POSIX only */
void pthread_setconcurrency_np(int con);
int pthread_getconcurrency_np(void);
void pthread_yield_np(void);

pthread_attr_t detached_attr;
pthread_mutex_t output_print_lk;
pthread_mutex_t global_count_lk;

int          global_count = 0;

work_t       *work_q = NULL;
pthread_cond_t work_q_cv;
pthread_mutex_t work_q_lk;
pthread_mutex_t debug_lock;

#include "debug.h" /* must be included AFTER the
                  mutex_t debug_lock line */

work_t       *search_q = NULL;
pthread_mutex_t search_q_lk;
pthread_cond_t search_q_cv;
int          search_pool_cnt = 0; /* the count in the pool now */
int          search_thr_limit = 0; /* the max in the pool */

work_t       *cascade_q = NULL;
pthread_mutex_t cascade_q_lk;
pthread_cond_t cascade_q_cv;
int          cascade_pool_cnt = 0;
int          cascade_thr_limit = 0;

int          running = 0;
pthread_mutex_t running_lk;

pthread_mutex_t stat_lk;
time_t       st_start = 0;
int          st_dir_search = 0;
int          st_file_search = 0;
int          st_line_search = 0;
int          st_cascade = 0;
int          st_cascade_pool = 0;
int          st_cascade_destroy = 0;
int          st_search = 0;
int          st_pool = 0;
int          st_maxrun = 0;
int          st_worknull = 0;
int          st_workfds = 0;
int          st_worklimit = 0;
int          st_destroy = 0;

int          all_done = 0;
int          work_cnt = 0;
```

EXAMPLE A-1 Source Code for `tgrep` Program (Continued)

```
int          current_open_files = 0;
int          tglimit = UNLIMITED; /* if -B limit the number of
                                threads */
int          progress_offset = 1;
int          progress = 0; /* protected by the print_lock ! */
unsigned int flags = 0;
int          regexp_cnt = 0;
char         *string[MAXREGEXP];
int          debug = 0;
int          use_pmatch = 0;
char         file_pat[255]; /* file patten match */
PATTERN     *pm_file_pat; /* compiled file target string (pmatch()) */

/*
 * Main: This is where the fun starts
 */
int
main(int argc, char **argv)
{
    int          c,out_thr_flags;
    long         max_open_files = 0l, ncpus = 0l;
    extern int   optind;
    extern char *optarg;
    int          prio = 0;
    struct stat sbuf;
    pthread_t    tid,dtid;
    void         *status;
    char         *e = NULL, *d = NULL; /* for debug flags */
    int          debug_file = 0;
    struct sigaction sigact;
    sigset_t     set,oset;
    int          err = 0, i = 0, pm_file_len = 0;
    work_t       *work;
    int          restart_cnt = 10;

    /* NO OTHER THREADS ARE RUNNING */
    flags = FR_RECUR; /* the default */

    while ((c = getopt(argc, argv, "d:e:bchilnsvwruf:p:BCSZzHP:")) != EOF) {
        switch (c) {
#ifdef DEBUG
            case 'd':
                debug = atoi(optarg);
                if (debug == 0)
                    debug_usage();

                d = optarg;
                fprintf(stderr,"tgrep: Debug on at level(s) ");
                while (*d) {
                    for (i=0; i<9; i++)
                        if (debug_set[i].level == *d) {
                            debug_levels |= debug_set[i].flag;
                            fprintf(stderr,"%c ",debug_set[i].level);
                        }
                }
                break;
#endif
        }
    }
}
```

EXAMPLE A-1 Source Code for `tgrep` Program (Continued)

```
        break;
    }
    d++;
}
fprintf(stderr, "\n");
break;
case 'f': debug_file = atoi(optarg); break;
#endif /* DEBUG */

case 'B':
    flags |= TG_BATCH;
#endifdef __lock_lint
/* locklint complains here, but there are no other threads */
    if ((e = getenv("TGLIMIT")) {
        tglimit = atoi(e);
    }
    else {
        if (!(flags & FS_NOERROR)) /* order dependent! */
            fprintf(stderr, "env TGLIMIT not set, overriding -B\n");
        flags &= ~TG_BATCH;
    }
#endif

    break;
case 'p':
    flags |= TG_FILEPAT;
    strcpy(file_pat, optarg);
    pm_file_pat = makepat(file_pat, NULL);
    break;
case 'P':
    flags |= TG_PROGRESS;
    progress_offset = atoi(optarg);
    break;
case 'S': flags |= FS_STATS;    break;
case 'b': flags |= FB_BLOCK;    break;
case 'c': flags |= FC_COUNT;    break;
case 'h': flags |= FH_HOLDNAME; break;
case 'i': flags |= FI_IGNCASE;  break;
case 'l': flags |= FL_NAMEONLY; break;
case 'n': flags |= FN_NUMBER;   break;
case 's': flags |= FS_NOERROR;  break;
case 'v': flags |= FV_REVERSE;  break;
case 'w': flags |= FW_WORD;     break;
case 'r': flags &= ~FR_RECUR;   break;
case 'C': flags |= FC_LINE;     break;
case 'e':
    if (regexp_cnt == MAXREGEXP) {
        fprintf(stderr, "Max number of regexp's (%d) exceeded!\n",
            MAXREGEXP);
        exit(1);
    }
    flags |= FE_REGEXP;
    if ((string[regexp_cnt] = (char *)malloc(strlen(optarg)+1)) == NULL) {
        fprintf(stderr, "tgrep: No space for search string(s)\n");
    }
}
```

EXAMPLE A-1 Source Code for `tgrep` Program (Continued)

```
        exit(1);
    }
    memset(string[regex_cnt], 0, strlen(optarg)+1);
    strcpy(string[regex_cnt], optarg);
    regex_cnt++;
    break;
case 'z':
case 'Z': regex_usage();
    break;
case 'H':
case '?':
default : usage();
    }
}
if (flags & FS_STATS)
    st_start = time(NULL);

if (!(flags & FE_REGEX)) {
    if (argc - optind < 1) {
        fprintf(stderr, "tgrep: Must supply a search string(s) "
            "and file list or directory\n");
        usage();
    }
    if ((string[0]=(char *)malloc(strlen(argv[optind])+1))==NULL) {
        fprintf(stderr, "tgrep: No space for search string(s)\n");
        exit(1);
    }
    memset(string[0], 0, strlen(argv[optind])+1);
    strcpy(string[0], argv[optind]);
    regex_cnt=1;
    optind++;
}

if (flags & FI_IGNCASE)
    for (i=0; i<regex_cnt; i++)
        uncase(string[i]);

if (flags & FE_REGEX) {
    for (i=0; i<regex_cnt; i++)
        pm_pat[i] = makepat(string[i], NULL);
    use_pmatch = 1;
}
else {
    bm_pat = bm_makepat(string[0]); /* only one allowed */
}

flags |= FX_STDIN;

max_open_files = sysconf(_SC_OPEN_MAX);
ncpus = sysconf(_SC_NPROCESSORS_ONLN);
if ((max_open_files - HOLD_FDS - debug_file) < 1) {
    fprintf(stderr, "tgrep: You MUST have at least ONE fd "
```

EXAMPLE A-1 Source Code for `tgrep` Program (Continued)

```
        "that can be used, check limit (>10)\n");
    exit(1);
}
search_thr_limit = max_open_files - HOLD_FDS - debug_file;
cascade_thr_limit = search_thr_limit / 2;
/* the number of files that can be open */
current_open_files = search_thr_limit;

pthread_attr_init(&detached_attr);
pthread_attr_setdetachstate(&detached_attr,
    PTHREAD_CREATE_DETACHED);

pthread_mutex_init(&global_count_lk,NULL);
pthread_mutex_init(&output_print_lk,NULL);
pthread_mutex_init(&work_q_lk,NULL);
pthread_mutex_init(&running_lk,NULL);
pthread_cond_init(&work_q_cv,NULL);
pthread_mutex_init(&search_q_lk,NULL);
pthread_cond_init(&search_q_cv,NULL);
pthread_mutex_init(&cascade_q_lk,NULL);
pthread_cond_init(&cascade_q_cv,NULL);

if ((argc == optind) && ((flags & TG_FILEPAT) || (flags & FR_RECUR))) {
    add_work(".",DIRT);
    flags = (flags & ~FX_STDIN);
}
for ( ; optind < argc; optind++) {
    restart_cnt = 10;
    flags = (flags & ~FX_STDIN);
    STAT_AGAIN:
    if (stat(argv[optind], &sbuf)) {
        if (errno == EINTR) { /* try again !, restart */
            if (--restart_cnt)
                goto STAT_AGAIN;
        }
        if (!(flags & FS_NOERROR))
            fprintf(stderr,"tgrep: Can't stat file/dir %s, %s\n",
                argv[optind], strerror(errno));
        continue;
    }
    switch (sbuf.st_mode & S_IFMT) {
    case S_IFREG :
        if (flags & TG_FILEPAT) {
            if (pmatch(pm_file_pat, argv[optind], &pm_file_len))
                DP(DLEVEL1,("File pat match %s\n",argv[optind]));
            add_work(argv[optind],FILET);
        }
        else {
            add_work(argv[optind],FILET);
        }
        break;
    case S_IFDIR :
        if (flags & FR_RECUR) {
```

EXAMPLE A-1 Source Code for `tgrep` Program (Continued)

```
        add_work(argv[optind],DIRT);
    }
    else {
        if (!(flags & FS_NOERROR))
            fprintf(stderr,"tgrep: Can't search directory %s, "
                    "-r option is on. Directory ignored.\n",
                    argv[optind]);
    }
    break;
}
}

pthread_setconcurrency_np(3);

if (flags & FX_STDIN) {
    fprintf(stderr,"tgrep: stdin option is not coded at this time\n");
    exit(0);
    /* XXX Need to fix this SOON */
    search_thr(NULL);
    if (flags & FC_COUNT) {
        pthread_mutex_lock(&global_count_lk);
        printf("%d\n",global_count);
        pthread_mutex_unlock(&global_count_lk);
    }
    if (flags & FS_STATS)
        prnt_stats();
    exit(0);
}

pthread_mutex_lock(&work_q_lk);
if (!work_q) {
    if (!(flags & FS_NOERROR))
        fprintf(stderr,"tgrep: No files to search.\n");
    exit(0);
}
pthread_mutex_unlock(&work_q_lk);

DP(DLEVEL1,("Starting to loop through the work_q for work\n"));

/* OTHER THREADS ARE RUNNING */
while (1) {
    pthread_mutex_lock(&work_q_lk);
    while ((work_q == NULL || current_open_files == 0 || tglimit <= 0) &&
            all_done == 0) {
        if (flags & FS_STATS) {
            pthread_mutex_lock(&stat_lk);
            if (work_q == NULL)
                st_worknull++;
            if (current_open_files == 0)
                st_workfds++;
            if (tglimit <= 0)
                st_worklimit++;
            pthread_mutex_unlock(&stat_lk);
        }
    }
}
```

EXAMPLE A-1 Source Code for `tgrep` Program (Continued)

```
pthread_cond_wait(&work_q_cv, &work_q_lk);
}
if (all_done != 0) {
    pthread_mutex_unlock(&work_q_lk);
    DP(DLEVEL1, ("All_done was set to TRUE\n"));
    goto OUT;
}
work = work_q;
work_q = work->next; /* maybe NULL */
work->next = NULL;
current_open_files--;
pthread_mutex_unlock(&work_q_lk);

tid = 0;
switch (work->tp) {
case DIRT:
    pthread_mutex_lock(&cascade_q_lk);
    if (cascade_pool_cnt) {
        if (flags & FS_STATS) {
            pthread_mutex_lock(&stat_lk);
            st_cascade_pool++;
            pthread_mutex_unlock(&stat_lk);
        }
        work->next = cascade_q;
        cascade_q = work;
        pthread_cond_signal(&cascade_q_cv);
        pthread_mutex_unlock(&cascade_q_lk);
        DP(DLEVEL2, ("Sent work to cascade pool thread\n"));
    }
    else {
        pthread_mutex_unlock(&cascade_q_lk);
        err = pthread_create(&tid, &detached_attr, cascade, (void *)work);
        DP(DLEVEL2, ("Sent work to new cascade thread\n"));
        if (flags & FS_STATS) {
            pthread_mutex_lock(&stat_lk);
            st_cascade++;
            pthread_mutex_unlock(&stat_lk);
        }
    }
    break;
case FILET:
    pthread_mutex_lock(&search_q_lk);
    if (search_pool_cnt) {
        if (flags & FS_STATS) {
            pthread_mutex_lock(&stat_lk);
            st_pool++;
            pthread_mutex_unlock(&stat_lk);
        }
        work->next = search_q; /* could be null */
        search_q = work;
        pthread_cond_signal(&search_q_cv);
        pthread_mutex_unlock(&search_q_lk);
        DP(DLEVEL2, ("Sent work to search pool thread\n"));
    }
}
```

EXAMPLE A-1 Source Code for tgrep Program (Continued)

```
    }
    else {
        pthread_mutex_unlock(&search_q_lk);
        err = pthread_create(&tid,&detached_attr,
                            search_thr, (void *)work);
        pthread_setconcurrency_np(pthread_getconcurrency_np()+1);
        DP(DLEVEL2, ("Sent work to new search thread\n"));
        if (flags & FS_STATS) {
            pthread_mutex_lock(&stat_lk);
            st_search++;
            pthread_mutex_unlock(&stat_lk);
        }
    }
    break;
default:
    fprintf(stderr, "tgrep: Internal error, work_t->tp not valid\n");
    exit(1);
}
if (err) { /* NEED TO FIX THIS CODE. Exiting is just wrong */
    fprintf(stderr, "Could not create new thread!\n");
    exit(1);
}
}

OUT:
if (flags & TG_PROGRESS) {
    if (progress)
        fprintf(stderr, ".\n");
    else
        fprintf(stderr, "\n");
}
/* we are done, print the stuff. All other threads are parked */
if (flags & FC_COUNT) {
    pthread_mutex_lock(&global_count_lk);
    printf("%d\n", global_count);
    pthread_mutex_unlock(&global_count_lk);
}
if (flags & FS_STATS)
    prnt_stats();
return(0); /* should have a return from main */
}

/*
 * Add_Work: Called from the main thread, and cascade threads to add file
 * and directory names to the work Q.
 */
int
add_work(char *path, int tp)
{
    work_t      *wt, *ww, *wp;

    if ((wt = (work_t *)malloc(sizeof(work_t))) == NULL)
        goto ERROR;
```

EXAMPLE A-1 Source Code for `tgrep` Program (Continued)

```
if ((wt->path = (char *)malloc(strlen(path)+1)) == NULL)
    goto ERROR;

strcpy(wt->path,path);
wt->tp = tp;
wt->next = NULL;
if (flags & FS_STATS) {
    pthread_mutex_lock(&stat_lk);
    if (wt->tp == DIRT)
        st_dir_search++;
    else
        st_file_search++;
    pthread_mutex_unlock(&stat_lk);
}
pthread_mutex_lock(&work_q_lk);
work_cnt++;
wt->next = work_q;
work_q = wt;
pthread_cond_signal(&work_q_cv);
pthread_mutex_unlock(&work_q_lk);
return(0);
ERROR:
if (!(flags & FS_NOERROR))
    fprintf(stderr,"tgrep: Could not add %s to work queue. Ignored\n",
           path);
return(-1);
}

/*
 * Search thread: Started by the main thread when a file name is found
 * on the work Q to be searched. If all the needed resources are ready
 * a new search thread will be created.
 */
void *
search_thr(void *arg) /* work_t *arg */
{
    FILE          *fin;
    char          fin_buf[(BUFSIZ*4)]; /* 4 Kbytes */
    work_t        *wt,std;
    int           line_count;
    char          rline[128];
    char          cline[128];
    char          *line;
    register char *p,*pp;
    int           pm_len;
    int           len = 0;
    long          byte_count;
    long          next_line;
    int           show_line; /* for the -v option */
    register int  slen,plen,i;
    out_t         *out = NULL; /* this threads output list */

    pthread_yield_np();
```

EXAMPLE A-1 Source Code for `tgrep` Program (Continued)

```
wt = (work_t *)arg; /* first pass, wt is passed to use. */

/* len = strlen(string);*/ /* only set on first pass */

while (1) { /* reuse the search threads */
    /* init all back to zero */
    line_count = 0;
    byte_count = 0l;
    next_line = 0l;
    show_line = 0;

    pthread_mutex_lock(&running_lk);
    running++;
    pthread_mutex_unlock(&running_lk);
    pthread_mutex_lock(&work_q_lk);
    tglimit--;
    pthread_mutex_unlock(&work_q_lk);
    DP(DLEVEL5, ("searching file (STDIO) %s\n", wt->path));

    if ((fin = fopen(wt->path, "r")) == NULL) {
        if (!(flags & FS_NOERROR)) {
            fprintf(stderr, "tgrep: %s. File \"%s\" not searched.\n",
                strerror(errno), wt->path);
        }
        goto ERROR;
    }
    setvbuf(fin, fin_buf, _IOFBF, (BUFSIZ*4)); /* XXX */
    DP(DLEVEL5, ("Search thread has opened file %s\n", wt->path));
    while ((fgets(rline, 127, fin)) != NULL) {
        if (flags & FS_STATS) {
            pthread_mutex_lock(&stat_lk);
            st_line_search++;
            pthread_mutex_unlock(&stat_lk);
        }
        slen = strlen(rline);
        next_line += slen;
        line_count++;
        if (rline[slen-1] == '\n')
            rline[slen-1] = '\0';
        /*
        ** If the uncase flag is set, copy the read in line (rline)
        ** To the uncase line (cline) Set the line pointer to point at
        ** cline.
        ** If the case flag is NOT set, then point line at rline.
        ** line is what is compared, rline is what is printed on a
        ** match.
        */
        if (flags & FI_IGNCASE) {
            strcpy(cline, rline);
            uncase(cline);
            line = cline;
        }
        else {
```

EXAMPLE A-1 Source Code for `tgrep` Program (Continued)

```
        line = rline;
    }
    show_line = 1; /* assume no match, if -v set */
    /* The old code removed */
    if (use_pmatch) {
        for (i=0; i<regexp_cnt; i++) {
            if (pmatch(pm_pat[i], line, &pm_len)) {
                if (!(flags & FV_REVERSE)) {
                    add_output_local(&out,wt,line_count,
                                    byte_count,rline);
                    continue_line(rline,fin,out,wt,
                                &line_count,&byte_count);
                }
                else {
                    show_line = 0;
                } /* end of if -v flag if / else block */
                /*
                ** if we get here on ANY of the regexp targets
                ** jump out of the loop, we found a single
                ** match so do not keep looking!
                ** If name only, do not keep searching the same
                ** file, we found a single match, so close the file,
                ** print the file name and move on to the next file.
                */
                if (flags & FL_NAMEONLY)
                    goto OUT_OF_LOOP;
                else
                    goto OUT_AND_DONE;
            } /* end found a match if block */
        } /* end of the for pat[s] loop */
    }
    else {
        if (bm_pmatch( bm_pat, line)) {
            if (!(flags & FV_REVERSE)) {
                add_output_local(&out,wt,line_count,byte_count,rline);
                continue_line(rline,fin,out,wt,
                            &line_count,&byte_count);
            }
            else {
                show_line = 0;
            }
            if (flags & FL_NAMEONLY)
                goto OUT_OF_LOOP;
        }
    }
    OUT_AND_DONE:
    if ((flags & FV_REVERSE) && show_line) {
        add_output_local(&out,wt,line_count,byte_count,rline);
        show_line = 0;
    }
    byte_count = next_line;
}
OUT_OF_LOOP:
```

EXAMPLE A-1 Source Code for tgrep Program (Continued)

```
fclose(fin);
/*
** The search part is done, but before we give back the FD,
** and park this thread in the search thread pool, print the
** local output we have gathered.
*/
print_local_output(out,wt); /* this also frees out nodes */
out = NULL; /* for the next time around, if there is one */
ERROR:
DP(DLEVEL5, ("Search done for %s\n", wt->path));
free(wt->path);
free(wt);

notrun();
pthread_mutex_lock(&search_q_lk);
if (search_pool_cnt > search_thr_limit) {
    pthread_mutex_unlock(&search_q_lk);
    DP(DLEVEL5, ("Search thread exiting\n"));
    if (flags & FS_STATS) {
        pthread_mutex_lock(&stat_lk);
        st_destroy++;
        pthread_mutex_unlock(&stat_lk);
    }
    return(0);
}
else {
    search_pool_cnt++;
    while (!search_q)
        pthread_cond_wait(&search_q_cv, &search_q_lk);
    search_pool_cnt--;
    wt = search_q; /* we have work to do! */
    if (search_q->next)
        search_q = search_q->next;
    else
        search_q = NULL;
    pthread_mutex_unlock(&search_q_lk);
}
}
/*NOTREACHED*/
}

/*
* Continue line: Special case search with the -C flag set. If you are
* searching files like Makefiles, some lines might have escape char's to
* continue the line on the next line. So the target string can be found, but
* no data is displayed. This function continues to print the escaped line
* until there are no more "\" chars found.
*/
int
continue_line(char *rline, FILE *fin, out_t *out, work_t *wt,
              int *lc, long *bc)
{
    int len;
```

EXAMPLE A-1 Source Code for `tgrep` Program (Continued)

```
int cnt = 0;
char *line;
char nline[128];

if (!(flags & FC_LINE))
    return(0);

line = rline;
AGAIN:
len = strlen(line);
if (line[len-1] == '\\') {
    if ((fgets(nline,127,fin)) == NULL) {
        return(cnt);
    }
    line = nline;
    len = strlen(line);
    if (line[len-1] == '\\n')
        line[len-1] = '\\0';
    *bc = *bc + len;
    *lc++;
    add_output_local(&out,wt,*lc,*bc,line);
    cnt++;
    goto AGAIN;
}
return(cnt);
}

/*
 * cascade: This thread is started by the main thread when directory names
 * are found on the work Q. The thread reads all the new file, and directory
 * names from the directory it was started when and adds the names to the
 * work Q. (it finds more work!)
 */

void *
cascade(void *arg) /* work_t *arg */
{
    char    fullpath[1025];
    int     restart_cnt = 10;
    DIR     *dp;

    char    dir_buf[sizeof(struct dirent) + PATH_MAX];
    struct dirent *dent = (struct dirent *)dir_buf;
    struct stat sbuf;
    char    *fpath;
    work_t  *wt;
    int     fl = 0, dl = 0;
    int     pm_file_len = 0;

    pthread_yield_np(); /* try to give control back to main thread */
    wt = (work_t *)arg;

    while(1) {
```

EXAMPLE A-1 Source Code for `tgrep` Program (Continued)

```
fl = 0;
dl = 0;
restart_cnt = 10;
pm_file_len = 0;

pthread_mutex_lock(&running_lk);
running++;
pthread_mutex_unlock(&running_lk);
pthread_mutex_lock(&work_q_lk);
tglimit--;
pthread_mutex_unlock(&work_q_lk);

if (!wt) {
    if (!(flags & FS_NOERROR))
        fprintf(stderr, "tgrep: Bad work node passed to cascade\n");
    goto DONE;
}
fpath = (char *)wt->path;
if (!fpath) {
    if (!(flags & FS_NOERROR))
        fprintf(stderr, "tgrep: Bad path name passed to cascade\n");
    goto DONE;
}
DP(DLEVEL3, ("Cascading on %s\n", fpath));
if ((dp = opendir(fpath)) == NULL) {
    if (!(flags & FS_NOERROR))
        fprintf(stderr, "tgrep: Can't open dir %s, %s. Ignored.\n",
                fpath, strerror(errno));
    goto DONE;
}
while ((readdir_r(dp, dent)) != NULL) {
    restart_cnt = 10; /* only try to restart the interrupted 10 X */

    if (dent->d_name[0] == '.') {
        if (dent->d_name[1] == '.' && dent->d_name[2] == '\0')
            continue;
        if (dent->d_name[1] == '\0')
            continue;
    }

    fl = strlen(fpath);
    dl = strlen(dent->d_name);
    if ((fl + 1 + dl) > 1024) {
        fprintf(stderr, "tgrep: Path %s/%s is too long. "
                "MaxPath = 1024\n",
                fpath, dent->d_name);
        continue; /* try the next name in this directory */
    }
    strcpy(fullpath, fpath);
    strcat(fullpath, "/");
    strcat(fullpath, dent->d_name);

    RESTART_STAT:
```

EXAMPLE A-1 Source Code for `tgrep` Program (Continued)

```
if (stat(fullpath,&sbuf)) {
    if (errno == EINTR) {
        if (--restart_cnt)
            goto RESTART_STAT;
    }
    if (!(flags & FS_NOERROR))
        fprintf(stderr,"tgrep: Can't stat file/dir %s, %s. "
            "Ignored.\n",
            fullpath,strerror(errno));
    goto ERROR;
}

switch (sbuf.st_mode & S_IFMT) {
case S_IFREG :
    if (flags & TG_FILEPAT) {
        if (pmatch(pm_file_pat, dent->d_name, &pm_file_len)) {
            DP(DLEVEL3,("file pat match (cascade) %s\n",
                dent->d_name));
            add_work(fullpath,FILET);
        }
    }
    else {
        add_work(fullpath,FILET);
        DP(DLEVEL3,("cascade added file (MATCH) %s to Work Q\n",
            fullpath));
    }
    break;

case S_IFDIR :
    DP(DLEVEL3,("cascade added dir %s to Work Q\n",fullpath));
    add_work(fullpath,DIRT);
    break;
}
}

ERROR:
    closedir(dp);

DONE:
    free(wt->path);
    free(wt);
    notrun();
    pthread_mutex_lock(&cascade_q_lk);
    if (cascade_pool_cnt > cascade_thr_limit) {
        pthread_mutex_unlock(&cascade_q_lk);
        DP(DLEVEL5,("Cascade thread exiting\n"));
        if (flags & FS_STATS) {
            pthread_mutex_lock(&stat_lk);
            st_cascade_destroy++;
            pthread_mutex_unlock(&stat_lk);
        }
        return(0); /* pthread_exit */
    }
}
```

EXAMPLE A-1 Source Code for `tgrep` Program (Continued)

```
else {
    DP(DLEVEL5, ("Cascade thread waiting in pool\n"));
    cascade_pool_cnt++;
    while (!cascade_q)
        pthread_cond_wait(&cascade_q_cv, &cascade_q_lk);
    cascade_pool_cnt--;
    wt = cascade_q; /* we have work to do! */
    if (cascade_q->next)
        cascade_q = cascade_q->next;
    else
        cascade_q = NULL;
    pthread_mutex_unlock(&cascade_q_lk);
}
}
/*NOTREACHED*/
}

/*
 * Print Local Output: Called by the search thread after it is done searching
 * a single file. If any output was saved (matching lines), the lines are
 * displayed as a group on stdout.
 */
int
print_local_output(out_t *out, work_t *wt)
{
    out_t      *pp, *op;
    int        out_count = 0;
    int        printed = 0;

    pp = out;
    pthread_mutex_lock(&output_print_lk);
    if (pp && (flags & TG_PROGRESS)) {
        progress++;
        if (progress >= progress_offset) {
            progress = 0;
            fprintf(stderr, ".");
        }
    }
    while (pp) {
        out_count++;
        if (!(flags & FC_COUNT)) {
            if (flags & FL_NAMEONLY) { /* Print name ONLY ! */
                if (!printed) {
                    printed = 1;
                    printf("%s\n", wt->path);
                }
            }
            else { /* We are printing more than just the name */
                if (!(flags & FH_HOLDNAME))
                    printf("%s :", wt->path);
                if (flags & FB_BLOCK)
                    printf("%ld:", pp->byte_count/512+1);
                if (flags & FN_NUMBER)
```

EXAMPLE A-1 Source Code for `tgrep` Program (Continued)

```
        printf("%d:",pp->line_count);
        printf("%s\n",pp->line);
    }
}
op = pp;
pp = pp->next;
/* free the nodes as we go down the list */
free(op->line);
free(op);
}

pthread_mutex_unlock(&output_print_lk);
pthread_mutex_lock(&global_count_lk);
global_count += out_count;
pthread_mutex_unlock(&global_count_lk);
return(0);
}

/*
 * add output local: is called by a search thread as it finds matching lines.
 * the matching line, its byte offset, line count, etc. are stored until the
 * search thread is done searching the file, then the lines are printed as
 * a group. This way the lines from more than a single file are not mixed
 * together.
 */

int
add_output_local(out_t **out, work_t *wt,int lc, long bc, char *line)
{
    out_t      *ot,*oo, *op;

    if (( ot = (out_t *)malloc(sizeof(out_t))) == NULL)
        goto ERROR;
    if (( ot->line = (char *)malloc(strlen(line)+1)) == NULL)
        goto ERROR;

    strcpy(ot->line,line);
    ot->line_count = lc;
    ot->byte_count = bc;

    if (!*out) {
        *out = ot;
        ot->next = NULL;
        return(0);
    }
    /* append to the END of the list; keep things sorted! */
    op = oo = *out;
    while(oo) {
        op = oo;
        oo = oo->next;
    }
    op->next = ot;
    ot->next = NULL;
}
```

EXAMPLE A-1 Source Code for tgrep Program (Continued)

```
return(0);

ERROR:
    if (!(flags & FS_NOERROR))
        fprintf(stderr,"tgrep: Output lost. No space. "
            "[%s: line %d byte %d match : %s\n",
            wt->path,lc,bc,line);
    return(1);
}

/*
 * print stats: If the -S flag is set, after ALL files have been searched,
 * main thread calls this function to print the stats it keeps on how the
 * search went.
 */

void
prnt_stats(void)
{
    float a,b,c;
    float t = 0.0;
    time_t st_end = 0;
    char    tl[80];

    st_end = time(NULL); /* stop the clock */
    printf("\n----- Tgrep Stats. ----- \n");
    printf("Number of directories searched:      %d\n",st_dir_search);
    printf("Number of files searched:                %d\n",st_file_search);
    c = (float)(st_dir_search + st_file_search) / (float)(st_end - st_start);
    printf("Dir/files per second:                    %3.2f\n",c);
    printf("Number of lines searched:                 %d\n",st_line_search);
    printf("Number of matching lines to target:       %d\n",global_count);

    printf("Number of cascade threads created:         %d\n",st_cascade);
    printf("Number of cascade threads from pool:      %d\n",st_cascade_pool);
    a = st_cascade_pool; b = st_dir_search;
    printf("Cascade thread pool hit rate:             %3.2f%%\n",((a/b)*100));
    printf("Cascade pool overall size:                %d\n",cascade_pool_cnt);
    printf("Number of search threads created:         %d\n",st_search);
    printf("Number of search threads from pool:       %d\n",st_pool);
    a = st_pool; b = st_file_search;
    printf("Search thread pool hit rate:              %3.2f%%\n",((a/b)*100));
    printf("Search pool overall size:                 %d\n",search_pool_cnt);
    printf("Search pool size limit:                   %d\n",search_thr_limit);
    printf("Number of search threads destroyed:       %d\n",st_destroy);

    printf("Max # of threads running concurrently:    %d\n",st_maxrun);
    printf("Total run time, in seconds.                %d\n",
        (st_end - st_start));

    /* Why did we wait ? */
    a = st_workfds; b = st_dir_search+st_file_search;
    c = (a/b)*100; t += c;
}
```

EXAMPLE A-1 Source Code for `tgrep` Program (Continued)

```
printf("Work stopped due to no FD's:  (%.3d)          %d Times, %3.2f%%\n",
      search_thr_limit,st_workfds,c);
a = st_worknull; b = st_dir_search+st_file_search;
c = (a/b)*100; t += c;
printf("Work stopped due to no work on Q:          %d Times, %3.2f%%\n",
      st_worknull,c);
if (tlimit == UNLIMITED)
    strcpy(tl,"Unlimited");
else
    sprintf(tl,"  %.3d  ",tlimit);
a = st_worklimit; b = st_dir_search+st_file_search;
c = (a/b)*100; t += c;
printf("Work stopped due to TGLIMIT:  (%.9s) %d Times, %3.2f%%\n",
      tl,st_worklimit,c);
printf("Work continued to be handed out:          %3.2f%%\n",100.00-t);
printf("-----\n");
}
/*
 * not running: A glue function to track if any search threads or cascade
 * threads are running. When the count is zero, and the work Q is NULL,
 * we can safely say, WE ARE DONE.
 */
void
notrun (void)
{
    pthread_mutex_lock(&work_q_lk);
    work_cnt--;
    tlimit++;
    current_open_files++;
    pthread_mutex_lock(&running_lk);
    if (flags & FS_STATS) {
        pthread_mutex_lock(&stat_lk);
        if (running > st_maxrun) {
            st_maxrun = running;
            DP(DLEVEL6,("Max Running has increased to %d\n",st_maxrun));
        }
        pthread_mutex_unlock(&stat_lk);
    }
    running--;
    if (work_cnt == 0 && running == 0) {
        all_done = 1;
        DP(DLEVEL6,("Setting ALL_DONE flag to TRUE.\n"));
    }
    pthread_mutex_unlock(&running_lk);
    pthread_cond_signal(&work_q_cv);
    pthread_mutex_unlock(&work_q_lk);
}
/*
 * uncase: A glue function. If the -i (case insensitive) flag is set, the
 * target string and the read in line is converted to lower case before
 * comparing them.
 */
```


EXAMPLE A-1 Source Code for `tgrep` Program (Continued)

```
fprintf(stderr, "    if you supply a directory name.\n");
fprintf(stderr, "    If you do not supply a file, or directory name,\n");
fprintf(stderr, "    and the -r option is not set, the current\n");
fprintf(stderr, "    directory \".\" will be used.\n");
fprintf(stderr, "    All the other options should work \"like\" grep\n");
fprintf(stderr, "    The -p patten is regexp; tgrep will search only\n");
fprintf(stderr, "\n");
fprintf(stderr, "    Copy Right By Ron Winacott, 1993-1995.\n");
fprintf(stderr, "\n");
exit(0);
}

/*
 * regexp usage: Tell the world about tgrep custom (THREAD SAFE) regexp!
 */
int
regexp_usage (void)
{
    fprintf(stderr, "usage: tgrep <options> -e \"pattern\" <-e ...> "
           "<{file,dir}>...\n");
    fprintf(stderr, "\n");
    fprintf(stderr, "metachars:\n");
    fprintf(stderr, "    . - match any character\n");
    fprintf(stderr, "    * - match 0 or more occurrences of previous char\n");
    fprintf(stderr, "    + - match 1 or more occurrences of previous char.\n");
    fprintf(stderr, "    ^ - match at beginning of string\n");
    fprintf(stderr, "    $ - match end of string\n");
    fprintf(stderr, "    [ - start of character class\n");
    fprintf(stderr, "    ] - end of character class\n");
    fprintf(stderr, "    ( - start of a new pattern\n");
    fprintf(stderr, "    ) - end of a new pattern\n");
    fprintf(stderr, "    @(n)c - match <c> at column <n>\n");
    fprintf(stderr, "    | - match either pattern\n");
    fprintf(stderr, "    \\ - escape any special characters\n");
    fprintf(stderr, "    \\c - escape any special characters\n");
    fprintf(stderr, "    \\o - turn on any special characters\n");
    fprintf(stderr, "\n");
    fprintf(stderr, "To match two different patterns in the same command\n");
    fprintf(stderr, "Use the or function. \n"
           "ie: tgrep -e \"(pat1)|(pat2)\" file\n"
           "This will match any line with \"pat1\" or \"pat2\" in it.\n");
    fprintf(stderr, "You can also use up to %d -e expressions\n", MAXREGEXP);
    fprintf(stderr, "RegExp Pattern matching brought to you by Marc Staveley\n");
    exit(0);
}

/*
 * debug usage: If compiled with -DDEBUG, turn it on, and tell the world
 * how to get tgrep to print debug info on different threads.
 */

#ifdef DEBUG
void
```

EXAMPLE A-1 Source Code for `tgrep` Program (Continued)

```
debug_usage(void)
{
    int i = 0;

    fprintf(stderr, "DEBUG usage and levels:\n");
    fprintf(stderr, "-----\n");
    fprintf(stderr, "Level          code\n");
    fprintf(stderr, "-----\n");
    fprintf(stderr, "0          This message.\n");
    for (i=0; i<9; i++) {
        fprintf(stderr, "%d          %s\n", i+1, debug_set[i].name);
    }
    fprintf(stderr, "-----\n");
    fprintf(stderr, "You can or the levels together like -d134 for levels\n");
    fprintf(stderr, "1 and 3 and 4.\n");
    fprintf(stderr, "\n");
    exit(0);
}
#endif

/* Pthreads NP functions */

#ifdef __sun
void
pthread_setconcurrency_np(int con)
{
    thr_setconcurrency(con);
}

int
pthread_getconcurrency_np(void)
{
    return(thr_getconcurrency());
}

void
pthread_yield_np(void)
{
    /*      In Solaris 2.4, these functions always return - 1 and set errno to ENOSYS */
    if (sched_yield()) /* call UI interface if we are older than 2.5 */
        thr_yield();
}

#else
void
pthread_setconcurrency_np(int con)
{
    return;
}

int
pthread_getconcurrency_np(void)
{

```

EXAMPLE A-1 Source Code for tgrep Program *(Continued)*

```
        return(0);
    }

void
pthread_yield_np(void)
{
    return;
}
#endif
```

Solaris Threads Example: barrier.c

The `barrier.c` program demonstrates an implementation of a barrier for Solaris threads. (See “Parallelizing a Loop on a Shared-Memory Parallel Computer” on page 223 for a definition of barriers.)

EXAMPLE B-1 Solaris Threads Example: barrier.c

```
#define _REENTRANT

/* Include Files      */

#include <thread.h>
#include <errno.h>

/* Constants & Macros  */

/* Data Declarations  */

typedef struct {
    int    maxcnt;      /* maximum number of runners */
    struct _sb {
        cond_t wait_cv; /* cv for waiters at barrier */
        mutex_t wait_lk; /* mutex for waiters at barrier */
        int    runners; /* number of running threads */
    } sb[2];
    struct _sb *sbp;    /* current sub-barrier */
} barrier_t;

/*
 * barrier_init - initialize a barrier variable.
 *
 */

int
barrier_init( barrier_t *bp, int count, int type, void *arg ) {
    int n;
```

EXAMPLE B-1 Solaris Threads Example: barrier.c (Continued)

```
int i;

if (count < 1)
    return(EINVAL);

bp->maxcnt = count;
bp->sbp = &bp->sb[0];

    for (i = 0; i < 2; ++i) {
#if defined(__cplusplus)
    struct barrier_t::_sb *sbp = &( bp->sb[i] );
#else
    struct _sb *sbp = &( bp->sb[i] );
#endif
    sbp->runners = count;

    if (n = mutex_init(&sbp->wait_lk, type, arg))
        return(n);

    if (n = cond_init(&sbp->wait_cv, type, arg))
        return(n);
    }
    return(0);
}

/*
 * barrier_wait - wait at a barrier for everyone to arrive.
 *
 */

int
barrier_wait(register barrier_t *bp) {
#if defined(__cplusplus)
    register struct barrier_t::_sb *sbp = bp->sbp;
#else
    register struct _sb *sbp = bp->sbp;
#endif
    mutex_lock(&sbp->wait_lk);

    if (sbp->runners == 1) { /* last thread to reach barrier */
        if (bp->maxcnt != 1) {
            /* reset runner count and switch sub-barriers */
            sbp->runners = bp->maxcnt;
            bp->sbp = (bp->sbp == &bp->sb[0])
                ? &bp->sb[1] : &bp->sb[0];

            /* wake up the waiters */
            cond_broadcast(&sbp->wait_cv);
        }
    } else {
        sbp->runners--; /* one less runner */

        while (sbp->runners != bp->maxcnt)
```

EXAMPLE B-1 Solaris Threads Example: barrier.c (Continued)

```
        cond_wait( &sbp->wait_cv, &sbp->wait_lk);
    }

    mutex_unlock(&sbp->wait_lk);

    return(0);
}

/*
 * barrier_destroy - destroy a barrier variable.
 *
 */

int
barrier_destroy(barrier_t *bp) {
    int    n;
    int    i;

    for (i=0; i < 2; ++ i) {
        if (n = cond_destroy(&bp->sb[i].wait_cv))
            return( n );

        if (n = mutex_destroy( &bp->sb[i].wait_lk))
            return(n);
    }

    return(0);
}

#define NTHR    4
#define NCOMPUTATION 2
#define NITER    1000
#define NSQRT    1000

void *
compute(barrier_t *ba )
{
    int count = NCOMPUTATION;

    while (count--) {
        barrier_wait( ba );
        /* do parallel computation */
    }
}

main( int argc, char *argv[] ) {
    int    i;
    int    niter;
    int    nthr;
    barrier_t    ba;
    double    et;
    thread_t    *tid;
```

EXAMPLE B-1 Solaris Threads Example: barrier.c *(Continued)*

```
switch ( argc ) {
    default:
        case 3 :      niter   = atoi( argv[1] );
                     nthr    = atoi( argv[2] );
                     break;

        case 2 :      niter   = atoi( argv[1] );
                     nthr    = NTHR;
                     break;

        case 1 :      niter   = NITER;
                     nthr    = NTHR;
                     break;
}

barrier_init( &ba, nthr + 1, USYNC_THREAD, NULL );
tid = (thread_t *) calloc(nthr, sizeof(thread_t));

for (i = 0; i < nthr; ++i) {
    int    n;

    if (n = thr_create(NULL, 0,
        (void (*)( void *)) compute,
        &ba, NULL, &tid[i])) {
        errno = n;
        perror("thr_create");
        exit(1);
    }
}

for (i = 0; i < NCOMPUTATION; i++) {
    barrier_wait(&ba );
    /* do parallel algorithm */
}

for (i = 0; i < nthr; i++) {
    thr_join(tid[i], NULL, NULL);
}
}
```

Index

Numbers and Symbols

- 32-bit architectures, 70
- 64-bit environment
 - data type model, 23
 - /dev/kmem, 23
 - /dev/mem, 23
 - large file support, 24
 - large virtual address space, 23
 - libkvm, 23
 - libraries, 23
 - /proc restrictions, 23
 - registers, 24

A

- Ada, 151
- adding
 - signals to mask, 41
- aio_errno, 156
- AIO_INPROGRESS, 156
- aio_result_t, 156
- aiocancel(3AIO), 155
- aioread(3AIO), 155
- aiowait(3AIO), 156
- aiowrite(3AIO), 155
- algorithms
 - faster with MT, 17
 - parallel, 224
 - sequential, 224
- ANSI C, 171
- application-level threads, 16

- architecture
 - multiprocessor, 220
 - SPARC, 223
 - SPARC, 70, 221
- assert statement, 113, 214
- Async-Signal-Safe
 - signal handlers, 151
 - functions, 149, 162
- asynchronous
 - event notification, 116
 - I/O, 154
 - semaphore use, 116
 - signals, 145, 149
- atomic, defined, 70
- automatic
 - stack allocation, 66

B

- binary semaphores, 115
- binding
 - threads to LWPs, 189
 - values to keys, 30, 194
- bottlenecks, 216
- bound threads, 16, 20
 - defined, 16
 - reasons to bind, 22

C

- C++, 171

- cache, defined, 220
- caching
 - threads data structure, 217
- changing the signal mask, 40, 192
- coarse-grained locking, 213
- code lock, 212
- code monitor, 212, 214
- compile flag
 - D_POSIX_C_SOURCE, 167
 - D_POSIX_PTHREAD_SEMANTICS, 167
 - D_REENTRANT, 167
 - single-threaded application, 167
- completion semantics, 150
- cond_broadcast(3THR), 200, 202
- cond_destroy(3THR), 200
- cond_init(3THR), 199, 205
- cond_signal(3THR), 200
- cond_timedwait(3THR), 201
- cond_wait(3THR), 154, 200
- condition variables, 70, 98, 114
- condition wait
 - POSIX threads, 153
 - Solaris threads, 153
- contention, 215
- continue execution, 181
- counting semaphores, 16, 115
- creating
 - stacks, 66, 189, 191
 - thread-specific keys, 30, 194
 - threads, 26, 28, 217
- critical section, 222
- custom stack, 66, 190

D

- daemon threads, 190
- data
 - global, 30
 - local, 30
 - lock, 212
 - races, 159
 - shared, 20, 222
 - thread specific
 - See thread-specific data
- dbx(1), 171
- deadlock, 214
- debugging, 169, 173

- debugging (*continued*)
 - asynchronous signals, 170
 - dbx(1), 171
 - deadlocks, 170
 - hidden gap in synchronization, 170
 - inadequate stack size, 170
 - large automatic arrays, 170
 - long-jumping without releasing mutex lock, 170
 - mdb(1), 171
 - no synchronization of global memory, 169
 - passing pointer to caller's stack, 169
 - recursive deadlock, 170
 - reevaluate conditions after return from condition wait, 170
- deleting signals from mask, 41
- destructor function, 31, 36
- detached threads, 28, 53, 189
- Dijkstra, E. W., 114

E

- EAGAIN, 27, 31, 90, 92, 104, 120, 190
- EBUSY, 88, 92, 104, 129, 184, 186
- EDEADLK, 28, 90
- EFAULT, 88, 183
- EINTR, 119, 139, 145, 154
- EINVAL, 27, 30, 32, 37, 40, 45, 52, 57, 59, 67, 73, 78, 85, 88, 93, 100, 104, 109, 117, 119, 123, 131, 183, 190
- ENOMEM, 31, 72, 91, 93, 100, 104, 190
- ENOSPC, 117
- ENOSYS, 38, 78, 85, 89
- ENOTRECOVERABLE, 91, 93
- ENOTSUP, 39, 59, 61, 78, 85
- EOWNERDEAD, 90, 92
- EPERM, 78, 91, 117
- errno, 34, 167, 169, 210
- __errno, 169
- errno
 - global variables, 210
- error checking, 40
- ESRCH, 28, 30, 39, 45, 180
- ETIME, 108
- event notification, 116
- examining the signal mask, 40, 192
- exec(2), 136, 138

exit(2), 139,190
exit(3C), 42

F

fair share scheduler (FSS) scheduling class, 144
finding
 minimum stack size, 191
 thread priority, 195
fine-grained locking, 213
fixed priority scheduling class (FX), 144
flags to thr_create(), 189
flockfile(3C), 157
flowchart of compile options, 168
fork(2), 137,139,200
fork1(2), 138
FORTRAN, 171
funlockfile(3C), 157

G

getc(3C), 157
getc_unlocked(3C), 157
gethostbyname(3NSL), 210
gethostbyname_r(3NSL), 211
getrusage(3C), 142
global
 data, 213
 side effects, 216
 state, 212
 variables, 33,209

H

heap, malloc(3C) storage from, 28

I

I/O
 asynchronous, 154
 nonsequential, 156
 standard, 157
 synchronous, 154
inheriting priority, 188

interrupt, 145
invariants, 113,213

J

joining threads, 27,53,193

K

keys
 bind value to key, 194
 get specific key, 33,194
 global into private, 35
 storing value of, 33,194
kill(2), 145,147

L

/lib/libc, 163,165,168
/lib/libC, 163
/lib/libdl_stubs, 163
/lib/libintl, 163,165
/lib/libm, 163,165
/lib/libmalloc, 163,165
/lib/libmapmalloc, 163,165
/lib/libnsl, 163,165,169
/lib/libpthread, 165,168
/lib/libresolv, 163
/lib/librt, 165
/lib/libsocket, 163,165
/lib/libthread, 19,165,168,217
/lib/libw, 163,165
/lib/libX11, 163
/lib/strtoaddr, 163
libraries
 MT-Safe, 163
library
 C routines, 209
 threads, 165,217
lightweight processes, 20,142,217
 creation, 219
 debugging, 171
 defined, 16
 independence, 218
 not supported, 20

- lightweight processes (*continued*)
 - special capabilities, 218
- limits, resources, 142
- linking with libpthread
 - lc, 167
 - ld, 167
 - lpthread, 167
- linking with libthread
 - lc, 167
 - ld, 167
 - lthread, 167
- ln(1)
 - linking, 165
- local variable, 211
- lock hierarchy, 215
- locking, 212
 - coarse grained, 213, 216
 - code, 212
 - conditional, 95
 - data, 213
 - fine-grained, 213, 216
 - guidelines, 215
 - invariants, 213
- locks, 70
 - mutual exclusion, 70, 97, 137
 - read-write, 187
 - readers/writer, 70
- longjmp(3C), 142, 251
- lposix4 library
 - POSIX 1003.1c semaphore, 168
- lseek(2), 157
- LWPs, *See* lightweight processes

M

- main(), 216
- malloc(3C), 28
- MAP_NORESERVE, 65
- MAP_SHARED, 139
- mdb(1), 171
- memory
 - consistency, 219
 - ordering, relaxed, 222
 - strongly ordered, 221
- mmap(2), 65, 139
- monitor, code, 212, 214
- mprotect(2), 191

- mt, 168
- MT-Safe libraries
 - alternative mmap(2)-based memory
 - allocation library, 163
 - C++ runtime shared objects, 163
 - internationalization, 163
 - math library, 163
 - network interfaces of the form
 - getXXbyYY_r, 163
 - network name-to-address translation, 163
 - socket library for making network
 - connections, 163
 - space-efficient memory allocation, 163
 - static switch compiling, 163
 - thread-safe form of unsafe interfaces, 163
 - thread-specific errno support, 163
 - wide character and wide string support for
 - multibyte locales, 163
 - X11 Windows routines, 163
- multiple-readers, single-writer locks, 187
- multiprocessors, 219, 224
- multithreading
 - defined, 16
- mutex
 - PTHREAD_MUTEX_ERRORCHECK, 89
 - PTHREAD_MUTEX_NORMAL, 89
 - PTHREAD_MUTEX_RECURSIVE, 89
- mutex, mutual exclusion locks, 214
- mutex_destroy(3THR), 197
- mutex_init(3THR), 196, 205
- mutex_lock(3THR), 198
- mutex scope, 73
- mutex_trylock(3THR), 198, 215
- mutex_unlock(3THR), 198
- mutual exclusion locks, 70, 97, 137
 - attributes, 72
 - deadlock, 94
 - default attributes, 70
 - scope, Solaris and POSIX, 71
 - type attribute, 75

N

- NDEBUG, 113
- netdir, 163
- netselect, 163
- nice(2), 143

nondetached threads, 42
 nonsequential I/O, 156
 null
 threads, 66, 190
 null procedures
 /lib/libpthread stub, 168
 /lib/libthread stub, 168

P

parallel
 algorithms, 224
 Pascal, 171
 PC
 program counter, 19
 PC_GETCID, 143
 PC_GETCLINFO, 143
 PC_GETPARMS, 143
 PC_SETPARMS, 143
 Peterson's Algorithm, 222
 PL/1 language, 146
 portability, 70
 POSIX 1003.4a, 17
 pread(2), 155
 printf(3S), 151
 printf problem, 211
 priocntl(2), 143
 PC_GETCID, 143
 PC_GETCLINFO, 143
 PC_GETPARMS, 143
 PC_SETPARMS, 143
 priority, 19, 143
 finding for a thread, 195
 inheritance, 188, 195
 and scheduling, 195
 range, 195
 setting for a thread, 195
 priority inversion, 77
 process
 terminating, 42
 traditional UNIX, 15
 producer/consumer problem, 131, 206, 221
 programmer-allocated stack, 66, 191
 prolagen
 semaphore, P operation, 115
 pthread_atfork(3THR), 41, 138
 pthread_attr_destroy(3THR), 52

pthread_attr_getdetachstate(3THR), 53
 pthread_attr_getguardsize(3THR), 55
 pthread_attr_getinheritsched(3THR), 61
 pthread_attr_getschedparam(3THR), 62
 pthread_attr_getschedpolicy(3THR), 59
 pthread_attr_getscope(3THR), 57
 pthread_attr_getstackaddr(3THR), 68
 pthread_attr_getstacksize(3THR), 65
 pthread_attr_init(3THR), 50
 attribute values, 51
 pthread_attr_setdetachstate(3THR), 52
 pthread_attr_setguardsize(3THR), 54
 pthread_attr_setinheritsched(3THR), 60
 pthread_attr_setschedparam(3THR), 61
 pthread_attr_setschedpolicy(3THR), 58
 pthread_attr_setscope(3THR), 56
 pthread_attr_setstackaddr(3THR), 67
 pthread_attr_setstacksize(3THR), 64
 pthread_cancel(3THR), 44
 pthread_cleanup_pop(3THR), 47
 pthread_cleanup_push(3THR), 47
 pthread_cond_broadcast(3THR), 104,
 109, 111, 145
 example, 110
 pthread_cond_destroy(3THR), 110
 pthread_cond_init(3THR), 103
 pthread_cond_signal(3THR), 104, 111,
 145
 example, 107
 pthread_cond_timedwait(3THR), 107
 example, 108
 pthread_cond_wait(3THR), 104, 111, 145
 example, 107
 pthread_condattr_destroy(3THR), 100
 pthread_condattr_getpshared(3THR), 102
 pthread_condattr_init(3THR), 99
 pthread_condattr_setpshared(3THR), 101
 pthread_create(3THR), 26
 pthread_detach(3THR), 29
 pthread_equal(3THR), 36
 pthread_exit(3THR), 41
 pthread_getconcurrency(3THR), 58
 pthread_getschedparam(3THR), 39
 pthread_getspecific(3THR), 33
 pthread_join(3THR), 27, 65, 155
 pthread_key_create(3THR), 30, 35
 example, 35
 pthread_key_delete(3THR), 31

pthread_kill(3THR), 39, 147
 pthread_mutex_consistent_np(3THR), 88
 pthread_mutex_destroy(3THR), 93
 pthread_mutex_getprioceiling(3THR)
 get priority ceiling of mutex, 83
 pthread_mutex_init(3THR), 87
 pthread_mutex_lock(3THR), 89
 example, 94, 96
 pthread_mutex_lock(3THR)
 example, 98
 pthread_mutex_setprioceiling(3THR)
 set priority ceiling of mutex, 82
 pthread_mutex_trylock(3THR), 92, 96
 pthread_mutex_unlock(3THR), 91
 example, 94, 97
 pthread_mutex_unlock(3THR)
 example, 98
 pthread_mutexattr_destroy(3THR), 72
 pthread_mutexattr_getprioceiling(3THR)
 get priority ceiling of mutex attribute, 80
 pthread_mutexattr_getprotocol(3THR)
 get protocol of mutex attribute, 79
 pthread_mutexattr_getpshared(3THR), 74
 pthread_mutexattr_getrobust_np(3THR)
 get robust attribute of mutex, 85
 pthread_mutexattr_gettype(3THR), 76
 pthread_mutexattr_init(3THR), 72
 pthread_mutexattr_setprioceiling(3THR)
 set priority ceiling of mutex attribute, 79
 pthread_mutexattr_setprotocol(3THR)
 set protocol of mutex attribute, 76
 pthread_mutexattr_setpshared(3THR), 73
 pthread_mutexattr_setrobust_np(3THR)
 set robust attribute of mutex, 83
 pthread_mutexattr_settype(3THR), 75
 pthread_once(3THR), 37
 PTHREAD_PRIO_INHERIT, 77
 PTHREAD_PRIO_NONE, 77
 PTHREAD_PRIO_PROTECT, 78
 pthread_rwlock_destroy(3THR), 130
 pthread_rwlock_init(3THR), 126
 pthread_rwlock_rdlock(3THR), 127
 pthread_rwlock_tryrdlock(3THR), 128
 pthread_rwlock_trywrlock(3THR), 129
 pthread_rwlock_unlock(3THR), 129
 pthread_rwlock_wrlock(3THR), 128
 pthread_rwlockattr_destroy(3THR), 123
 pthread_rwlockattr_getpshared(3THR), 125
 pthread_rwlockattr_init(3THR), 123
 pthread_rwlockattr_setpshared(3THR), 124
 PTHREAD_SCOPE_PROCESS, 21, 56
 PTHREAD_SCOPE_SYSTEM, 22, 56
 pthread_self(3THR), 36
 pthread_setcancelstate(3THR), 45
 pthread_setcanceltype(3THR), 45
 pthread_setconcurrency(3THR), 57
 pthread_setschedparam(3THR), 38
 pthread_setspecific(3THR), 32, 35
 example, 35
 pthread_sigmask(3THR), 40, 147
 PTHREAD_STACK_MIN(), 66
 pthread_testcancel(3THR), 46
 putc(3C), 157
 putc_unlocked(3C), 157
 pwrite(2), 155

R

_r, 211
 read(2), 156
 read-write locks, 125, 187
 attribute, 124
 attributes, 122
 readers/writer locks, 70
 realtime
 scheduling, 143
 red zone, 66, 191
 reentrant, 212
 described, 212
 functions, 161
 strategies for making, 212
 register state, 19
 relaxed memory ordering, 222
 remote procedure call RPC, 18
 replacing signal mask, 41
 resume execution, 181
 RPC, 18, 163, 217
 RT,, *See* realtime
 rw_rdlock(3THR), 183
 rw_tryrdlock(3THR), 184
 rw_trywrlock(3THR), 185
 rw_unlock(3THR), 186
 rw_wrlock(3THR), 185
 rwlock_destroy(3THR), 186
 rwlock_init(3THR), 182, 205

S

- SA_RESTART, 154
- safety, threads interfaces, 159, 164
- sched_yield(3RT), 37
- scheduling
 - realtime, 143
 - system class, 142
 - timeshare, 143
- scheduling class
 - fair share scheduler (FSS), 144
 - fixed priority scheduler (FX), 144
 - priority, 142
 - timeshare, 144
- sem_destroy(3RT), 120
- sem_init(3RT), 116
 - example, 121
- sem_post(3RT), 115, 118
 - example, 121
- sem_trywait(3RT), 115, 119
- sem_wait(3RT), 115, 119
 - example, 121
- sema_destroy(3THR), 204
- sema_init(3THR), 202, 205
- sema_post(3THR), 162
- sema_post(3THR), 203
- sema_trywait(3THR), 204
- sema_wait(3THR), 204
- semaphores, 70, 114, 133
 - binary, 115
 - counting, 115
 - counting, defined, 16
 - decrement semaphore value, 115
 - increment semaphore value, 115
 - interprocess, 117
 - intraprocess, 117
 - named, 118
- sending signal to thread, 39, 192
- sequential algorithms, 223
- setjmp(3C), 142, 150
- shared data, 20, 213
- shared-memory multiprocessor, 221
- SIG_BLOCK, 41
- SIG_DFL, 145
- SIG_IGN, 145
- SIG_SETMASK, 41
- SIG_UNBLOCK, 41
- sigaction(2), 145, 154
- SIGFPE, 146, 151
- SIGILL, 146
- SIGINT, 146, 150, 154
- SIGIO, 146, 156
- siglongjmp(3C), 151
- signal(3C), 145
- signal(5), 145
- signal.h, 40, 192
- signals
 - access mask, 40, 192
 - add to mask, 41
 - asynchronous, 145, 149
 - delete from mask, 41
 - handler, 145, 149
 - inheritance, 188
 - masks, 19
 - pending, 181, 188
 - replace current mask, 41
 - send to thread, 39, 192
 - SIG_BLOCK, 41
 - SIG_SETMASK, 41
 - SIG_UNBLOCK, 41
 - SIGSEGV, 65
 - unmasked and caught, 153
- sigprocmask(2), 147
- SIGPROF
 - interval timer, 140
- sigqueue(3RT), 145
- SIGSEGV, 65, 146
- sigsend(2), 145
- sigsetjmp(3C), 151
- sigtimedwait(3RT), 149
- SIGVTALRM
 - interval timer, 140
- sigwait(2), 148, 151
- single-threaded
 - assumptions, 209
 - code, 70
 - defined, 16
 - processes, 139
- size of stack, 64, 66, 189, 191
- stack, 217
 - address, 67, 189
 - boundaries, 65
 - creation, 67, 189
 - custom, 190
 - deallocation, 191
 - minimum size, 66, 191
 - overflows, 66

- stack (*continued*)
 - parameters, 28
 - pointer, 19
 - programmer-allocated, 66, 191
 - red zone, 66, 191
 - returning a pointer to, 161
 - size, 64, 66, 189, 191
- stack_base, 67, 189
- stack_size, 64, 189
- standard I/O, 157
- standards, 17
- start_routine(), 189
- static storage, 169, 209
- stdio, 34, 167
- store buffer, 223
- storing thread key value, 33, 194
- streaming a tape drive, 155
- strongly ordered memory, 221
- suspending a new thread, 189
- swap space, 65
- synchronization objects, 69, 133
 - condition variables, 70, 98, 114
 - mutex locks, 70, 97
 - read-write locks, 187
 - semaphores, 70, 114, 131, 202, 206
- synchronous I/O, 154
- system calls
 - handling errors, 210
- system scheduling class, 142

T

- __t_errno, 169
- tape drive, streaming, 155
- terminating
 - process, 42
 - threads, 28
- THR_BOUND, 189
- thr_continue(3THR), 181, 189
- thr_create(3THR), 188, 190
- THR_DAEMON, 190
- THR_DETACHED, 189
- thr_exit(3THR), 190, 192
- thr_getprio(3THR), 195
- thr_getspecific(3THR), 194
- thr_join(3THR), 193
- thr_keycreate(3THR), 194

- thr_kill(3THR), 162
- thr_kill(3THR), 192
- thr_min_stack(3THR), 189
- thr_self(3THR), 191
- thr_setprio(3THR), 195
- thr_setspecific(3THR), 194
- thr_sigsetmask(3THR), 162
- thr_sigsetmask(3THR), 192
- THR_SUSPENDED, 189
- thr_yield(3THR), 192, 215
- thread-directed signal, 149
- thread-private storage, 20
- thread-specific data, 30, 36
 - global, 33
 - global into private, 34
 - new storage class, 210
 - private, 33
- thread synchronization
 - condition variables, 23
 - mutex locks, 22
 - mutual exclusion locks, 70
 - read-write locks, 122
 - read/write locks, 22
 - semaphores, 23, 114
- threads
 - creating, 26, 28, 188, 190, 217
 - daemon, 190
 - defined, 16
 - detached, 28, 53, 189
 - exit status, 26
 - identifiers, 36, 42, 189
 - initial, 42
 - joining, 27, 42, 193
 - keys, 194
 - library, 165, 217
 - lightweight processes, 20
 - nondetached, 42
 - null, 66, 190
 - priority, 188
 - private data, 30
 - safety, 159, 164
 - signals, 153
 - stack, 161
 - suspended, 181
 - suspending, 189
 - synchronizing, 70, 133
 - terminating, 28, 41, 192
 - thread-specific data, 210

threads (*continued*)
 unbound threads, 21
 user-level, 16, 19
time-out, 201
 example, 108
timeshare scheduling class, 143
tiuser.h, 169
TLI, 163
TLI, 169
tools
 dbx(1), 171
 mdb(1), 171
total store order, 223
trap, 145
 default action, 146
TS,, *See* timeshare scheduling class
TSD, *See* thread-specific data

U

unbound threads, 142
 caching, 217
 priorities, 142
 and scheduling, 142
 prctl(2), 143
UNIX, 15, 17, 19, 145, 154, 156, 210
user-level threads, 16, 19
/usr/include/errno.h, 165
/usr/include/limits.h, 165
/usr/include/pthread.h, 165
/usr/include/signal.h, 165
/usr/include/thread.h, 165
/usr/include/unistd.h, 165
/usr/lib
 32-bit threads library, Solaris 9 Operating
 Environment, 169
/usr/lib/lwp
 32-bit threads library, Solaris 8 Operating
 Environment, 169
/usr/lib/lwp/64
 64-bit threads library, Solaris 8 Operating
 Environment, 169
USYNC_PROCESS, 182, 196, 199, 203, 205
USYNC_PROCESS_ROBUST, 196
USYNC_THREAD, 182, 196, 199, 203, 205

V

variables
 condition, 70, 98, 114, 133
 global, 209
 primitive, 70
verhogen
 semaphore, V operation, 115
vfork(2), 138

W

write(2), 156

X

XDR, 163

