



Linux[®] Resource Administration Guide

007-4413-016

COPYRIGHT

© 2002–2007, 2010, 2011, SGI. All rights reserved; provided portions may be copyright in third parties, as indicated elsewhere herein. No permission is granted to copy, distribute, or create derivative works from the contents of this electronic documentation in any manner, in whole or in part, without the prior written permission of SGI.

LIMITED RIGHTS LEGEND

The software described in this document is "commercial computer software" provided with restricted rights (except as to included open/free source) as specified in the FAR 52.227-19 and/or the DFAR 227.7202, or successive sections. Use beyond license provisions is a violation of worldwide intellectual property laws, treaties and conventions. This document is provided with limited rights as defined in 52.227-14.

TRADEMARKS AND ATTRIBUTIONS

SGI, the SGI logo, and Supportfolio are trademarks or registered trademarks of Silicon Graphics International Corp. or its subsidiaries in the United States and other countries.

Linux is a registered trademark of Linus Torvalds. UNIX and the X Window System are registered trademarks of The Open Group in the United States and other countries. All other trademarks mentioned herein are the property of their respective owners.

New Features in This Manual

This rewrite of the *Linux Resource Administration Guide* supports the SGI Performance Suite 1.3 release.

Major Documentation Changes

Added "Configuring Array Services for Use with MPI" on page 6.

Updated information in "Boot Cpuset" on page 76.

Updated "Configuring a User Cpuset for Interactive Sessions" on page 79.

Added Procedure 2-5 on page 79.

Record of Revision

Version	Description
001	February 2003 Original publication.
002	June 2003 Updated to support the SGI ProPack for Linux v2.2 release
003	October 2003 Updated to support the SGI ProPack for Linux v2.3 release.
004	February 2004 Updated to support the SGI ProPack for Linux v2.4 release.
005	May 2004 Updated to support the SGI ProPack 3 for Linux release.
006	January 2005 Updated to support the SGI ProPack 3 for Linux Service Pack 3 release.
007	February 2005 Updated to support the SGI ProPack 4 for Linux release.
008	April 2005 Updated to support the SGI ProPack 3 for Linux Service Pack 5 release.
009	August 2005 Updated to support the SGI ProPack 4 for Linux Service Pack 2 release.
010	January 2006 Updated to support the SGI ProPack 4 for Linux Service Pack 3 release.
011	July 2006 Updated to support the SGI ProPack 5 for Linux release.

Record of Revision

- 012 April 2007
Updated to support the SGI ProPack 5 for Linux Service Pack 1 release.
- 013 July 2008
Updated to support the SGI ProPack 6 for Linux release.
- 014 January 2009
Updated to support the SGI ProPack 6 for Linux Service Pack 2 release.
- 015 October 2010
Updated to support the SGI Performance Suite 1.0 release.
- 016 November 2011
Updated to support the SGI Performance Suite 1.3 release.

Contents

About This Guide	xvii
Related Publications	xvii
Obtaining Publications	xvii
Conventions	xvii
Reader Comments	xviii
1. Array Services	1
Finding the Array Services Release Notes	2
Array Services Package	3
Installing and Configuring Array Services for Single Host Systems	3
Installing and Configuring Array Services for Cluster or Partitioned Systems	4
Automatic Array Services Installation Steps	6
Configuring Array Services for Use with MPI	6
Using an Array	7
Using an Array System	8
Finding Basic Usage Information	8
Logging In to an Array	8
Invoking a Program	9
Managing Local Processes	11
Monitoring Local Processes and System Usage	11
Scheduling and Killing Local Processes	11
Summary of Local Process Management Commands	12
Using Array Services Commands	12

About Array Sessions	13
About Names of Arrays and Nodes	13
About Authentication Keys	14
Summary of Common Command Options	14
Specifying a Single Node	16
Common Environment Variables	16
Interrogating the Array	17
Learning Array Names	17
Learning Node Names	18
Learning Node Features	18
Learning User Names and Workload	19
Learning User Names	19
Learning Workload	19
Managing Distributed Processes	20
About Array Session Handles (ASH)	20
Listing Processes and ASH Values	21
Controlling Processes	22
Using arshell	22
About the Distributed Example	23
Managing Session Processes	24
About Job Container IDs	25
About Array Configuration	25
Security Considerations for Standard Array Services	26
About the Uses of the Configuration File	27
About Configuration File Format and Contents	28
Loading Configuration Data	28
About Substitution Syntax	29

- Testing Configuration Changes 30
- Configuring Arrays and Machines 31
 - Specifying Arrayname and Machine Names 31
 - Specifying IP Addresses and Ports 32
 - Specifying Additional Attributes 32
- Configuring Authentication Codes 32
- Configuring Array Commands 33
 - Operation of Array Commands 33
 - Summary of Command Definition Syntax 34
 - Configuring Local Options 36
 - Designing New Array Commands 37
- Secure Array Services 39
 - Differences between Standard and Secure Array Services 39
 - Secure Array Services Certificates 40
 - Secure Array Services Parameters 44
 - Secure Shell Considerations 44
- 2. Cpusets on Linux 47**
- An Overview of the Advantages Gained by Using Cpusets 47
 - Linux 2.6 Kernel Support for Cpusets 49
 - Cpuset Facility Capabilities 50
 - Initializing Cpusets 50
 - How to Determine if Cpusets are Installed 51
 - Fine-grained Control within Cpusets 52
 - Cpuset Interaction with Other Placement Mechanism 52
 - Cpusets and Thread Placement 54
 - Safe Job Migration and Cpusets 54

Cpuset File System Directories	56
Exclusive Cpusets	62
Notify on Release Flag	62
Memory Pressure of a Cpuset	63
Memory Spread	65
Memory Migration	66
Mask Format	67
List Format	67
Cpuset Permissions	67
CPU Scheduling and Memory Allocation for Cpusets	68
Linux Kernel CPU and Memory Placement Settings	68
Manipulating Cpusets	70
Using Cpusets at the Shell Prompt	70
Cpuset Command Line Utility	72
Boot Cpuset	76
Creating a Bootcpuset	77
bootcpuset.conf File	78
Configuring a User Cpuset for Interactive Sessions	79
Cpuset Text Format	82
Modifying the CPUs in a Cpuset and Kernel Processing	83
Using Cpusets with Hyper-Threads	84
Cpuset Programming Model	86
System Error Messages	88
3. NUMA Tools	91
Appendix A. Cpuset Library Functions	93
Extensible Application Programming Interface	93

Basic Cpuset Library Functions	94
cpuset_pin	95
cpuset_size	96
cpuset_where	96
cpuset_unpin	96
Advanced Cpuset Library Functions	97
cpuset_alloc	102
cpuset_free	103
cpuset_cpus_nbits	103
cpuset_mems_nbits	103
cpuset_setcpus	103
cpuset_setmems	104
cpuset_set_iopt	104
cpuset_set_sopt	105
cpuset_getcpus	105
cpuset_getmems	105
cpuset_cpus_weight	106
cpuset_mems_weight	106
cpuset_get_iopt	106
cpuset_get_sopt	107
cpuset_localcpus	107
cpuset_localmems	108
cpuset_cpumemdist	108
cpuset_cpu2node	108
cpuset_create	109
cpuset_delete	109
cpuset_query	109
cpuset_modify	110

Contents

cpuset_getcpusetpath	110
cpuset_cpusetofpid	111
cpuset_cpusetofpid	111
cpuset_mountpoint	112
cpuset_collides_exclusive	112
cpuset_init_pidlist	113
cpuset_pidlist_length	113
cpuset_get_pidlist	113
cpuset_free_pidlist	113
cpuset_move	114
cpuset_move_all	114
cpuset_move_cpuset_tasks	114
cpuset_migrate	115
cpuset_migrate_all	116
cpuset_reattach	116
cpuset_open_memory_pressure	116
cpuset_read_memory_pressure	117
cpuset_close_memory_pressure	117
cpuset_c_rel_to_sys_cpu	118
cpuset_c_sys_to_rel_cpu	118
cpuset_c_rel_to_sys_mem	118
cpuset_c_sys_to_rel_mem	118
cpuset_p_rel_to_sys_cpu	119
cpuset_p_sys_to_rel_cpu	119
cpuset_p_rel_to_sys_mem	119
cpuset_p_sys_to_rel_mem	119
cpuset_get_placement	119
cpuset_equal_placement	120

- cpuset_free_placement 120
- cpuset_cpupbind 120
- cpuset_latestcpu 121
- cpuset_membind 121
- cpuset_nuke 122
- cpuset_export 123
- cpuset_import 123
- cpuset_function 124
- cpuset_version 124
- Functions to Traverse a Cpuset Hierarchy 125**
 - cpuset_fts_open 125
 - cpuset_fts_read 125
 - cpuset_fts_reverse 126
 - cpuset_fts_rewind 126
 - cpuset_fts_get_path 126
 - cpuset_fts_get_stat 127
 - cpuset_fts_get_cpuset 127
 - cpuset_fts_get_errno 128
 - cpuset_fts_close 128
- Index 129**

Tables

Table 1-1	Information Sources for Invoking a Program	11
Table 1-2	Information Sources: Local Process Management	12
Table 1-3	Common Array Services Commands	13
Table 1-4	Array Services Environment Variables	17
Table 1-5	Information Sources: Array Configuration	26
Table 1-6	Subentries of a COMMAND Definition	34
Table 1-7	Substitutions Used in a COMMAND Definition	35
Table 1-8	Options of the COMMAND Definition	36
Table 1-9	Subentries of the LOCAL Entry	37

About This Guide

This guide is a reference document for people who manage the operation of SGI computer systems running SGI Performance Suite software. It contains information needed in the administration of various system resource management features.

This manual contains the following chapters:

- Chapter 1, "Array Services" on page 1
- Chapter 2, "Cpusets on Linux" on page 47
- Chapter 3, "NUMA Tools" on page 91
- Appendix A, "Cpuset Library Functions" on page 93

Related Publications

For a list of manuals supporting SGI Linux releases, see the *SGI Performance Suite 1.3 Start Here*.

Obtaining Publications

You can obtain SGI documentation in the following ways:

- See the SGI Technical Publications Library at: <http://docs.sgi.com>. Various formats are available. This library contains the most recent and most comprehensive set of online books, release notes, man pages, and other information.
- Online versions of the *SGI Performance Suite 1.3 Start Here* contains a list of currently active SGI software and hardware manuals.
- You can view man pages by typing `man title` on a command line.

Conventions

The following conventions are used throughout this document:

Convention	Meaning
<code>command</code>	This fixed-space font denotes literal items such as commands, files, routines, path names, signals, messages, and programming language structures.
<i>variable</i>	Italic typeface denotes variable entries and words or concepts being defined.
user input	This bold, fixed-space font denotes literal items that the user enters in interactive sessions. (Output is shown in nonbold, fixed-space font.)
[]	Brackets enclose optional portions of a command or directive line.
...	Ellipses indicate that a preceding element can be repeated.

Reader Comments

If you have comments about the technical accuracy, content, or organization of this publication, contact SGI. Be sure to include the title and document number of the publication with your comments. (Online, the document number is located in the front matter of the publication. In printed publications, the document number is located at the bottom of each page.)

You can contact SGI in any of the following ways:

- Send e-mail to the following address:
techpubs@sgi.com
- Contact your customer service representative and ask that an incident be filed in the SGI incident tracking system.
- Send mail to the following address:
SGI
Technical Publications
46600 Landing Parkway
Fremont, CA 94538

SGI values your comments and will respond to them promptly.

Array Services

Array Services includes administrator commands, libraries, daemons, and kernel extensions that support the execution of parallel applications across a number of hosts in a cluster, or *array*. The Message Passing Interface (MPI) of SGI® MPI uses Array Services to launch parallel applications. For information on MPI, see the *Message Passing Toolkit (MPT) User's Guide*.

The secure version of Array Services is built to make use of secure sockets layer (SSL) and secure shell (SSH).

Note: Differences between the standard version and the secure version of Array Services are noted throughout this chapter. For simplicity and clarity, the use of Array Services generally refers to both products. When noting differences between the two, a distinction is made between Array Services (AS), the standard product, and Secure Array Services (SAS), the security enhanced product.

The Array Services package requires that the process sets service be installed and running. This package is provided in the `sgi-procset` RPM.

When using SAS, you also need to install the `openssl` package available from the Linux distribution. For more information, see "Secure Array Services" on page 39.

Note: To install Secure Array Services, use the YaST Software Management and use the `Filter->search` function to search for secure array services by name (`sarraysvcs`).

A central concept in Array Services is the array session handle (ASH), a number that is used to logically group related processes that may be distributed across multiple systems. The ASH creates a global process namespace across the Array, facilitating accounting and administration

Array Services also provides an array configuration database, listing the nodes comprising an array. Array inventory inquiry functions provide a centralized, canonical view of the configuration of each node. Other array utilities let the administrator query and manipulate distributed array applications.

This chapter covers the following topics:

- "Finding the Array Services Release Notes" on page 2
- "Array Services Package" on page 3
- "Installing and Configuring Array Services for Single Host Systems" on page 3
- "Installing and Configuring Array Services for Cluster or Partitioned Systems" on page 4
- "Automatic Array Services Installation Steps" on page 6
- "Configuring Array Services for Use with MPI" on page 6
- "Using an Array" on page 7
- "Managing Local Processes" on page 11
- "Using Array Services Commands" on page 12
- "Summary of Common Command Options" on page 14
- "Interrogating the Array" on page 17
- "Managing Distributed Processes" on page 20
- "About Array Configuration" on page 25
- "Configuring Arrays and Machines" on page 31
- "Configuring Authentication Codes" on page 32
- "Configuring Array Commands" on page 33
- "Secure Array Services" on page 39

Finding the Array Services Release Notes

You can find information about the location of Array Services release notes in the description section of the RPM. For standard Array Services, enter the following:

```
sys:~ # rpm -ql sgi-arraysvcs | grep relnotes  
/usr/share/doc/packages/sgi-arraysvcs/README.relnotes
```

For Secure Array Services, enter the following:

```
sys:~ # rpm -ql sgi-sarraysvcs | grep relnotes
/usr/share/doc/packages/sgi-sarraysvcs/README.relnotes
```

Array Services Package

The Array Services package comprises the following primary components:

array daemon	Allocates ASH values and maintains information about node configuration and the relation of process IDs to ASHs. Array daemons reside on each node and work in cooperation.
array configuration database	Describes the array configuration used by array daemons and user programs. One copy at each node.
ainfo command	Lets the user or administrator query the Array configuration database and information about ASH values and processes.
array command	Executes a specified command on one or more nodes. Commands are predefined by the administrator in the configuration database.
arshell command	Starts a command remotely on a different node using the current ASH value.

The use of the `ainfo`, `array`, and `arshell` commands is covered in "Using an Array" on page 7.

Installing and Configuring Array Services for Single Host Systems

The normal SGI MPI system installation process installs and pre-configures Array Services software to enable single host Message Passing Toolkit (MPT) Message Passing Interface (MPI) jobs. The configuration steps encoded in the Array Services RPM installation script also automatically issue the `chkconfig(8)` commands that register the Array Services `arrayd(8)` daemon to be started upon system reboot. If the usual system reboot is done after installing a SGI MPI software release, you do not need to take any additional steps to configure Array Services.

Because there are two versions of the product, the standard version of Array Services is installed as `sgi-arraysvcs`. The security enhanced version is installed as `sgi-sarraysvc`. You cannot install both versions at the same time because they are mutually incompatible.

If you are installing a new Array Services RPM on a live system, the Array Services daemon should be stopped before upgrading the software and then restarted after the upgrade. To stop the standard Array Services daemon, perform the following command:

```
% /etc/init.d/array stop
```

To stop the secure Array Services daemon, perform the following command:

```
% /etc/init.d/sarray stop
```

To start the standard Array Services daemon without having to reboot your system, perform the following command:

```
% /etc/init.d/array start
```

To start the secure Array Services daemon without having to reboot your system, perform the following command:

```
% /etc/init.d/sarray start
```

The steps that are executed automatically by the Array Services RPM at install time are described in the Array Services release notes (for location of the release notes, see "Finding the Array Services Release Notes" on page 2) and in "Automatic Array Services Installation Steps" on page 6.

Installing and Configuring Array Services for Cluster or Partitioned Systems

On clustered or partitioned Altix systems, it is often desirable to enable MPT MPI jobs to execute on multiple hosts, rather than being confined to a single host.

Note: If you run secure Array Services, you also need to install the `openssl 0.9.7` package (or later) available from the Linux distribution. In addition, for the steps that follow, keep in mind that the daemon for SAS is named `sarrayd` versus `arrayd` on standard Array Services.

To configure Array Services to execute on multiple hosts, perform the following:

1. Identify a cluster name and a host list.

Edit the `/usr/lib/array/arrayd.conf` file to list the machines in your cluster. The `arrayd.conf` file allows many specifications. For information about these specifications, see the `arrayd.conf(4)` man page. The only required specifications that need to be configured are the name for the cluster and a list of hostnames in the cluster.

The `arrayd.conf.template` file can also be found under the `/usr/lib/array/` directory

In the following steps, changes are made to the `arrayd.conf` file so that the cluster is given the name `sgicluster` and it consists of hosts named `host1`, `host2`, and so on:

- a. Add an array entry that lists the host names one per line, as follows:

```
array sgicluster
    machine host1
    machine host2
    ....
```

- b. In the destination `array` directive, edit the default cluster name to be `sgicluster`, as follows:

```
destination array sgicluster
```

2. Choose an authentication policy: `NONE` or `SIMPLE`.

You need to choose the security level under which Array Services will operate. The choices are authentication settings of `NONE` or `SIMPLE`. Either way, start by commenting out the line in `/usr/lib/array/arrayd.auth` file that reads `AUTHENTICATION NOREMOTE`. If no authentication is required at your site, uncomment the `AUTHENTICATION NONE` line in the `arrayd.auth` file. If you choose simple authentication, create an `AUTHENTICATION SIMPLE` section as described in the `arrayd.auth(4)` man page.

Note: For sites concerned with security, `AUTHENTICATION SIMPLE` is a better choice. `AUTHENTICATION` is enforced to `NONE` when using secure Array Services and authentication is performed via certificate. For details, see "Secure Array Services" on page 39.

3. When you are configuring Secure Array Services, you need to configure certificate. For information how to do this, see "Secure Array Services Certificates" on page 40.

Automatic Array Services Installation Steps

The following steps are performed automatically during installation of the Array Services RPM:

- An account must exist on all hosts in the array for the purpose of running certain Array Services commands. This is controlled by the `/usr/lib/array/arrayd.conf` configuration file. The default is to use the user account `arraysvcs`. The account name can be changed in `arrayd.conf`. The user account `arraysvcs` is installed by default.
- The following entry is added to `/etc/services` file to define the `arrayd` service and port number. The default port number is 5434 and is specified in the `arrayd.conf` configuration file. Any value can be used for the port number, but all systems mentioned in the `arrayd.conf` file must use the same value.

```
sgi-arrayd  5434/tcp  # SGI Array Services daemon
```

- Standard Array Services are activated during installation with the `chkconfig(1)` command, as follows:

```
chkconfig --add array
```

Secure Array Services are activated during installation with the `chkconfig(1)` command, as follows:

```
chkconfig --add sarray
```

Configuring Array Services for Use with MPI

Array Services must be configured and running on your SGI system or cluster to perform the launch of MPI jobs. For SGI Altix ICE systems, the unique root file systems are best configured by the procedure outlined in the `arrayconfig_tempo(8)` man page. For other SGI systems and clusters, you can configure Array Services by executing the following command on one of the cluster hosts, making sure to list all host names on the command line:

```
/usr/sbin/arrayconfig -m host1 host2 ...
```


The `arrayconfig` command will configure `/etc/array/arrayd.conf` and `/etc/array/arrayd.auth` for your cluster. Copy these configuration files to every cluster host and then restart Array Services on every cluster host, as follows:

```
/etc/init.d/array restart
```

For a more elaborate configuration, consult the `arrayconfig(1)` and `arrayd.conf(4)` man pages and "Installing and Configuring Array Services for Cluster or Partitioned Systems" on page 4.

Using an Array

An Array system is an aggregation of *nodes*, that are servers bound together with a high-speed network and Array Services software. Array users have the advantage of greater performance and additional services. Array users access the system with familiar commands for job control, login and password management, and remote execution.

Array Services augments conventional facilities with additional services for array users and for array administrators. The extensions include support for global session management, array configuration management, batch processing, message passing, system administration, and performance visualization.

This section introduces the extensions for Array use, with pointers to more detailed information. The main topics are as follows:

- "Using an Array System" on page 8, summarizes what a user needs to know and the main facilities a user has available.
- "Managing Local Processes" on page 11, reviews the conventional tools for listing and controlling processes within one node.
- "Using Array Services Commands" on page 12, describes the common concepts, options, and environment variables used by the Array Services commands.
- "Interrogating the Array" on page 17, summarizes how to use Array Services commands to learn about the Array and its workload, with examples.
- "Summary of Common Command Options" on page 14
- "Managing Distributed Processes" on page 20, summarizes how to use Array Services commands to list and control processes in multiple nodes.

Using an Array System

The array system allows you to run distributed sessions on multiple nodes of an array. You can access the Array from either:

- A workstation
- An X terminal
- An ASCII terminal

In each case, you log in to one node of the Array in the way you would log in to any remote UNIX host. From a workstation or an X terminal you can of course open more than one terminal window and log into more than one node.

Finding Basic Usage Information

In order to use an Array, you need the following items of information:

- The name of the Array.

You use this *arrayname* in Array Services commands.

- The login name and password you will use on the Array.

You use these when logging in to the Array to use it.

- The hostnames of the array nodes.

Typically these names follow a simple pattern, often *arrayname1*, *arrayname2*, and so on.

- Any special resource-distribution or accounting rules that may apply to you or your group under a job scheduling system.

You can learn the hostnames of the array nodes if you know the array name, using the `ainfo` command as follows:

```
ainfo -a arrayname machines
```

Logging In to an Array

Each node in an Array has an associated hostname and IP network address. Typically, you use an Array by logging in to one node directly, or by logging in remotely from another host (such as the Array console or a networked workstation). For example,

from a workstation on the same network, this command would log you in to the node named `hydra6` as follows:

```
rlogin hydra6
```

For details of the `rlogin` command, see the `rlogin(1)` man page.

The system administrators of your array may choose to disallow direct node logins in order to schedule array resources. If your site is configured to disallow direct node logins, your administrators will be able to tell you how you are expected to submit work to the array—perhaps through remote execution software or batch queueing facilities.

Invoking a Program

Once you have access to an array, you can invoke programs of several classes:

- Ordinary (sequential) applications
- Parallel shared-memory applications within a node
- Parallel message-passing applications within a node
- Parallel message-passing applications distributed over multiple nodes (and possibly other servers on the same network running Array Services)

If you are allowed to do so, you can invoke programs explicitly from a logged-in shell command line; or you may use remote execution or a batch queueing system.

Programs that are X Windows clients must be started from an X server, either an X Terminal or a workstation running X Windows.

Some application classes may require input in the form of command line options, environment variables, or support files upon execution. For example:

- X client applications need the `DISPLAY` environment variable set to specify the X server (workstation or X-terminal) where their windows will display.
- A multithreaded program may require environment variables to be set describing the number of threads.

For example, C and Fortran programs that use parallel processing directives test the `MP_SET_NUMTHREADS` variable.

- Message Passing Interface (MPI) and Parallel Virtual Machine (PVM) message-passing programs may require support files to describe how many tasks to invoke on specified nodes.

Some information sources on program invocation are listed in Table 1-1 on page 11.

Table 1-1 Information Sources for Invoking a Program

Topic	Man Page
Remote login	rlogin(1)
Setting environment variables	environ(5), env(1)

Managing Local Processes

Each UNIX process has a *process identifier* (PID), a number that identifies that process within the node where it runs. It is important to realize that a PID is local to the node; so it is possible to have processes in different nodes using the same PID numbers.

Within a node, processes can be logically grouped in *process groups*. A process group is composed of a parent process together with all the processes that it creates. Each process group has a *process group identifier* (PGID). Like a PID, a PGID is defined locally to that node, and there is no guarantee of uniqueness across the Array.

Monitoring Local Processes and System Usage

You query the status of processes using the system command `ps`. To generate a full list of all processes on a local system, use a command such as the following:

```
ps -elfj
```

You can monitor the activity of processes using the command `top` (an ASCII display in a terminal window).

Scheduling and Killing Local Processes

You can schedule commands to run at specific times using the `at` command. You can kill or stop processes using the `kill` command. To destroy the process with PID 13032, use a command such as the following:

```
kill -KILL 13032
```

Summary of Local Process Management Commands

Table 1-2 on page 12, summarizes information about local process management..

Table 1-2 Information Sources: Local Process Management standard

Topic	Man Page
Process ID and process group	intro(2)
Listing and monitoring processes	ps(1), top(1)
Running programs at low priority	nice(1), batch(1)
Running programs at a scheduled time	at(1)
Terminating a process	kill(1)
Table for driving cron	crontab(5)
Schedule periodic background work	crontab(1P)

Using Array Services Commands

When an application starts processes on more than one node, the PID and PGID are no longer adequate to manage the application. The commands of Array Services give you the ability to view the entire array, and to control the processes of multinode programs.

Note: You can use Array Services commands from any workstation connected to an array system. You do not have to be logged in to an array node.

The following commands are common to Array Services operations as shown in Table 1-3 on page 13.

Note: The `arshell(1)` command is not installed or usable when you are running Secure Array Services.

Table 1-3 Common Array Services Commands

Topic	Man Page
Array Services Overview	array_services(5)
ainfo command	ainfo(1)
array command	Use array(1); configuration: arrayd.conf(4)
arshell command	arshell(1)
newsess command	newsess (1)

About Array Sessions

Array Services is composed of a daemon—a background process that is started at boot time in every node—and a set of commands such as `ainfo(1)`. The commands call on the daemon process in each node to get the information they need.

One concept that is basic to Array Services is the *array session*, which is a term for all the processes of one application, wherever they may execute. Normally, your login shell, with the programs you start from it, constitutes an array session. A batch job is an array session; and you can create a new shell with a new array session identity.

Each session is identified by an *array session handle* (ASH), a number that identifies any process that is part of that session. You use the ASH to query and to control all the processes of a program, even when they are running in different nodes.

About Names of Arrays and Nodes

Each node is server, and as such has a hostname. The hostname of a node is returned by the `hostname(1)` command executed in that node as follows:

```
% hostname
tokyo
```

The command is simple and documented in the `hostname(1)` man page. The more complicated issues of `hostname` syntax, and of how hostnames are resolved to hardware addresses are covered in `hostname(5)`.

An Array system as a whole has a name too. In most installations there is only a single Array, and you never need to specify which Array you mean. However, it is possible to have multiple Arrays available on a network, and you can direct Array Services commands to a specific Array.

About Authentication Keys

It is possible for the Array administrator to establish an authentication code, which is a 64-bit number, for all or some of the nodes in an array (see "Configuring Authentication Codes" on page 58). When this is done, each use of an Array Services command must specify the appropriate authentication key, as a command option, for the nodes it uses. Your system administrator will tell you if this is necessary.

Note: When running Secure Array Services, this configuration is not used. Authentication is enforced to `AUTHENTICATION_NONE`.

Summary of Common Command Options

The following Array Services commands have a consistent set of command options: `ainfo(1)`, `array(1)`, and `arshell(1)`. The default values of some options are set by environment variables listed in the next topic.

Note: The `arshell(1)` command is not installed or usable when you are running Secure Array Services.

For an `ainfo(1)` usage statement, perform the following:

```
sys:~ # ainfo
ainfo: Must specify a request
```

Usage: `ainfo <options> <request>`

Common Options:

```
-a <arrayname>  Send request to array <arrayname>
-b              Brief output (same as -f 2)
```


-h <ASH>	Array session handle
-i <PID>	Process identification number
-l	Displays local information only
-p <portnum>	Port address of daemon
-s <server>	Daemon host name or IP address
-t <timeout>	Time out (in seconds)
-f <format#>	Verbose mode

Valid requests (any case):

ARRAYS
ASH
DFLTARRAY
MACHINES
NEWASH
PIDS
PROCESSES
SERVER
SESSIONS

For an array(1) usage statement, perform the following:

```
sys:~ # array
```

Usage: array <options...> <command>

Common Options:

-a <arrayname>	Send request to array <arrayname>
-D	Send request directly to "-s" server
-F	Forward request to "-s" server via local server
-g	Run command in new array session
-Kl <key>	Local authentication key
-Kr <key>	Remote authentication key
-l	Local: do not broadcast to other machines
-p <portnum>	Connect to server on port <portnum>
-q	Quiet: do not print output from request
-s <server>	Connect directly to server on machine <server>
-t <timeout>	Time out request after <timeout> seconds

For an arshell(1) usage statement, perform the following:

```
sys:~ # arshell
```

```
Usage: arshell <options> [user@]host [-s] <command>
```

```
Common Options:
```

-D	Send request directly to server
-F	Forward request to server via local server
-Kl <key>	Local authentication key
-Kr <key>	Remote authentication key
-l <user>	Login as user <user>
-N	Do not revert to BSD rsh
-O	Force usage of BSD rsh
-X	Do not use ssh when available
-p <portnum>	Connect to server on port <portnum>
-s	Delimits end of arshell options/start of <command>

Specifying a Single Node

The `-l` and `-s` options work together. The `-l` (letter ell for “local”) option restricts the scope of a command to the node where the command is executed. By default, that is the node where the command is entered. When `-l` is not used, the scope of a query command is all nodes of the array. The `-s` (server, or node name) option directs the command to be executed on a specified node of the array. These options work together in query commands as follows:

- To interrogate all nodes as seen by the local node, use neither option.
- To interrogate only the local node, use only `-l`.
- To interrogate all nodes as seen by a specified node, use only `-s`.
- To interrogate only a particular node, use both `-s` and `-l`.

Common Environment Variables

The Array Services commands depend on environment variables to define default values for the less-common command options. These variables are summarized in Table 1-4.

Table 1-4 Array Services Environment Variables

Variable Name	Use	Default When Undefined
ARRAYD_FORWARD	When defined with a string starting with the letter <i>y</i> , all commands default to forwarding through the array daemon (option <i>-F</i>).	Commands default to direct communication (option <i>-D</i>).
ARRAYD_PORT	The port (socket) number monitored by the array daemon on the destination node.	The standard number of 5434, or the number given with option <i>-p</i> .
ARRAYD_LOCALKEY	Authentication key for the local node (option <i>-Kl</i>).	No authentication unless <i>-Kl</i> option is used.
ARRAYD_REMOTEKEY	Authentication key for the destination node (option <i>-Kr</i>).	No authentication unless <i>-Kr</i> option is used.
ARRAYD	The destination node, when not specified by the <i>-s</i> option.	The local node, or the node given with <i>-s</i> .

Interrogating the Array

Any user of an Array system can use Array Services commands to check the hardware components and the software workload of the Array. The commands needed are `ainfo` and `array`.

Learning Array Names

If your network includes more than one Array system, you can use `ainfo arrays` at one array node to list all the Array names that are configured, as in the following example.

```
homegrown% ainfo arrays
Arrays known to array services daemon
ARRAY DevArray
    IDENT 0x3381
ARRAY BigDevArray
    IDENT 0x7456
ARRAY test
```

```
IDENT 0x655e
```

Array names are configured into the array database by the administrator. Different Arrays might know different sets of other Array names.

Learning Node Names

You can use `ainfo machines` to learn the names and some features of all nodes in the current Array, as in the following example.

```
homegrown 175% ainfo -b machines
machine homegrown homegrown 5434 192.48.165.36 0
machine disarray disarray 5434 192.48.165.62 0
machine datarray datarray 5434 192.48.165.64 0
machine tokyo tokyo 5434 150.166.39.39 0
```

In this example, the `-b` option of `ainfo` is used to get a concise display.

Learning Node Features

You can use `ainfo nodeinfo` to request detailed information about one or all nodes in the array. To get information about the local node, use `ainfo -l nodeinfo`. However, to get information about only a particular other node, for example node `tokyo`, use `-l` and `-s`, as in the following example. (The example has been edited for brevity.)

```
homegrown 181% ainfo -s tokyo -l nodeinfo
Node information for server on machine "tokyo"
MACHINE tokyo
  VERSION 1.2
  8 PROCESSOR BOARDS
    BOARD: TYPE 15  SPEED 190
      CPU:  TYPE 9  REVISION 2.4
      FPU:  TYPE 9  REVISION 0.0
  ...
  16 IP INTERFACES  HOSTNAME tokyo  HOSTID 0xc01a5035
    DEVICE et0  NETWORK 150.166.39.0  ADDRESS 150.166.39.39  UP
    DEVICE atm0  NETWORK 255.255.255.255  ADDRESS 0.0.0.0  UP
    DEVICE atm1  NETWORK 255.255.255.255  ADDRESS 0.0.0.0  UP
  ...
  0 GRAPHICS INTERFACES
```

MEMORY

```
512 MB MAIN MEMORY
INTERLEAVE 4
```

If the `-l` option is omitted, the destination node will return information about every node that it knows.

Learning User Names and Workload

The system commands `who(1)`, `top(1)`, and `uptime(1)` are commonly used to get information about users and workload on one server. The `array(1)` command offers Array-wide equivalents to these commands.

Learning User Names

To get the names of all users logged in to the whole array, use `array who`. To learn the names of users logged in to a particular node, for example `tokyo`, use `-l` and `-s`, as in the following example. (The example has been edited for brevity and security.)

```
homegrown 180% array -s tokyo -l who
joecd    tokyo      frummage.eng.sgi -tcsh
joecd    tokyo      frummage.eng.sgi -tcsh
benf     tokyo      einstein.ued.sgi. /bin/tcsh
yohn     tokyo      rayleigh.eng.sg vi +153 fs/procfs/prd
...
```

Learning Workload

Two variants of the `array` command return workload information. The array-wide equivalent of `uptime` is `array uptime`, as follows:

```
homegrown 181% array uptime
homegrown: up 1 day, 7:40, 26 users, load average: 7.21, 6.35, 4.72
disarray:  up 2:53, 0 user, load average: 0.00, 0.00, 0.00
datarray:  up 5:34, 1 user, load average: 0.00, 0.00, 0.00
tokyo:     up 7 days, 9:11, 17 users, load average: 0.15, 0.31, 0.29
homegrown 182% array -l -s tokyo uptime
tokyo:     up 7 days, 9:11, 17 users, load average: 0.12, 0.30, 0.28
```

The command `array top` lists the processes that are currently using the most CPU time, with their ASH values, as in the following example.

```
homegrown 183% array top
      ASH      Host      PID User      %CPU Command
-----
0x1111ffff00000000 homegrown      5 root      1.20 vfs_sync
0x1111ffff000001e9 homegrown    1327 arraysvcs  1.19 atop
0x1111ffff000001e9 tokyo      19816 arraysvcs  0.73 atop
0x1111ffff000001e9 disarray     1106 arraysvcs  0.47 atop
0x1111ffff000001e9 datarray     1423 arraysvcs  0.42 atop
0x1111ffff00000000 homegrown      20 root      0.41 ShareII
0x1111ffff000000c0 homegrown   29683 kchang     0.37 ld
0x1111ffff0000001e homegrown    1324 root      0.17 arrayd
0x1111ffff00000000 homegrown     229 root      0.14 routed
0x1111ffff00000000 homegrown      19 root      0.09 pdflush
0x1111ffff000001e9 disarray     1105 arraysvcs  0.02 atopm
```

The `-l` and `-s` options can be used to select data about a single node, as usual.

Managing Distributed Processes

Using commands from Array Services, you can create and manage processes that are distributed across multiple nodes of the Array system.

About Array Session Handles (ASH)

In an Array system you can start a program with processes that are in more than one node. In order to name such collections of processes, Array Services software assigns each process to an *array session handle* (ASH).

An ASH is a number that is unique across the entire array (unlike a PID or PGID). An ASH is the same for every process that is part of a single array session; no matter which node the process runs in. You display and use ASH values with Array Services commands. Each time you log in to an Array node, your shell is given an ASH, which is used by all the processes you start from that shell.

The command `ainfo ash` returns the ASH of the current process on the local node, which is simply the ASH of the `ainfo` command itself.

```
homegrown 178% ainfo ash
Array session handle of process 10068: 0x1111ffff000002c1
homegrown 179% ainfo ash
```

```
Array session handle of process 10069: 0x1111ffff000002c1
```

In the preceding example, each instance of the `ainfo` command was a new process: first PID 10068, then PID 10069. However, the ASH is the same in both cases. This illustrates a very important rule: **every process inherits its parent's ASH**. In this case, each instance of `array` was forked by the command shell, and the ASH value shown is that of the shell, inherited by the child process.

You can create a new global ASH with the command `ainfo newash`, as follows:

```
homegrown 175% ainfo newash
Allocating new global ASH
0x11110000308b2f7c
```

This feature has little use at present. There is no existing command that can change its ASH, so you cannot assign the new ASH to another command. It is possible to write a program that takes an ASH from a command-line option and uses the Array Services function `setash()` to change to that ASH (however such a program must be privileged). No such program is distributed with Array Services.

Listing Processes and ASH Values

The command `array ps` returns a summary of all processes running on all nodes in an array. The display shows the ASH, the node, the PID, the associated username, the accumulated CPU time, and the command string.

To list all the processes on a particular node, use the `-l` and `-s` options. To list processes associated with a particular ASH, or a particular username, pipe the returned values through `grep`, as in the following example. (The display has been edited to save space.)

```
homegrown 182% array -l -s tokyo ps | fgrep wombat
0x261cffff0000054c      tokyo 19007  wombat    0:00 -csh
0x261cffff0000054a      tokyo 17940  wombat    0:00 csh -c (setenv...
0x261cffff0000054c      tokyo 18941  wombat    0:00 csh -c (setenv...
0x261cffff0000054a      tokyo 17957  wombat    0:44 xem -geometry 84x42
0x261cffff0000054a      tokyo 17938  wombat    0:00 rshd
0x261cffff0000054a      tokyo 18022  wombat    0:00 /bin/csh -i
0x261cffff0000054a      tokyo 17980  wombat    0:03 /usr/gnu/lib/ema...
0x261cffff0000054c      tokyo 18928  wombat    0:00 rshd
```

Controlling Processes

The `arshell` command lets you start an arbitrary program on a single other node. The `array` command gives you the ability to suspend, resume, or kill all processes associated with a specified ASH.

Using arshell

The `arshell` command is an Array Services extension of the familiar `rsh` command; it executes a single system command on a specified Array node. The difference from `rsh` is that the remote shell executes under the same ASH as the invoking shell (this is not true of simple `rsh`). The following example demonstrates the difference.

```
homegrown 179% ainfo ash
Array session handle of process 8506: 0x1111ffff00000425
homegrown 180% rsh arraysvcs@tokyo ainfo ash
Array session handle of process 13113: 0x261cffff0000145e
homegrown 181% arshell arraysvcs@tokyo ainfo ash
Array session handle of process 13119: 0x1111ffff00000425
```

You can use `arshell` to start a collection of unrelated programs in multiple nodes under a single ASH; then you can use the commands described under "Managing Session Processes" on page 24 to stop, resume, or kill them.

Both MPI and PVM use `arshell` to start up distributed processes.

Tip: The shell is a process under its own ASH. If you use the `array` command to stop or kill all processes started from a shell, you will stop or kill the shell also. In order to create a group of programs under a single ASH that can be killed safely, proceed as follows:

1. Within the new shell, start one or more programs using `arshell`.
2. Exit the nested shell.

Now you are back to the original shell. You know the ASH of all programs started from the nested shell. You can safely kill all jobs that have that ASH because the current shell is not affected.

About the Distributed Example

The programs launched with `arshell` are not coordinated (they could of course be written to communicate with each other, for example using sockets), and you must start each program individually.

The `array` command is designed to permit the simultaneous launch of programs on all nodes with a single command. However, `array` can only launch programs that have been configured into it, in the Array Services configuration file. (The creation and management of this file is discussed under "About Array Configuration" on page 25.)

In order to demonstrate process management in a simple way from the command line, the following command was inserted into the configuration file

`/usr/lib/array/arrayd.conf`:

```
#
# Local commands
#
command spin                # Do nothing on multiple machines
    invoke /usr/lib/array/spin
    user    %USER
    group   %GROUP
    options nowait
```

The invoked command, `/usr/lib/array/spin`, is a shell script that does nothing in a loop, as follows:

```
#!/bin/sh
# Go into a tight loop
#
interrupted() {
    echo "spin has been interrupted - goodbye"
    exit 0
}
trap interrupted 1 2
while [ ! -f /tmp/spin.stop ]; do
    sleep 5
done
echo "spin has been stopped - goodbye"
exit 1
```

With this preparation, the command `array spin` starts a process executing that script on every processor in the array. Alternatively, `array -l -s nodename spin` would start a process on one specific node.

Managing Session Processes

The following command sequence creates and then kills a `spin` process in every node. The first step creates a new session with its own ASH. This is so that later, `array kill` can be used without killing the interactive shell.

```
homegrown 175% ainfo ash
Array session handle of process 8912: 0x1111ffff0000032d
homegrown 175% ainfo ash
Array session handle of process 8941: 0x11110000308b2fa6
```

In the new session with ASH `0x11110000308b2fa6`, the command `array spin` starts the `/usr/lib/array/spin` script on every node. In this test array, there were only two nodes on this day, `homegrown` and `tokyo`.

```
homegrown 176% array spin
```

After exiting back to the original shell, the command `array ps` is used to search for all processes that have the ASH `0x11110000308b2fa6`.

```

homegrown 177% exit
homegrown 178% homegrown 177%
homegrown 177% ainfo ash
Array session handle of process 9257: 0x1111ffff0000032d
homegrown 179% array ps | fgrep 0x11110000308b2fa6
0x11110000308b2fa6 homegrown 9033 arraysvcs 0:00 /bin/sh /usr/lib/array/spin
0x11110000308b2fa6 homegrown 9618 arraysvcs 0:00 sleep 5
0x11110000308b2fa6      tokyo 26021 arraysvcs 0:00 /bin/sh /usr/lib/array/spin
0x11110000308b2fa6      tokyo 26072 arraysvcs 0:00 sleep 5
0x1111ffff0000032d homegrown 9642 arraysvcs 0:00 fgrep 0x11110000308b2fa6

```

There are two processes related to the `spin` script on each node. The next command kills them all.

```

homegrown 180% array kill 0x11110000308b2fa6
homegrown 181% array ps | fgrep 0x11110000308b2fa6
0x1111ffff0000032d homegrown 10030 arraysvcs 0:00 fgrep 0x11110000308b2fa6

```

The command `array suspend 0x11110000308b2fa6` would suspend the processes instead (however, it is hard to demonstrate that a `sleep` command has been suspended).

About Job Container IDs

Array systems have the capability to forward job IDs (JIDs) from the initiating host. All of the processes running in the ASH across one or more nodes in an array also belong to the same job.

When processes are running on the initiating host, they belong to the same job as the initiating process and operate under the limits established for that job. On remote nodes, a new job is created using the same JID as the initiating process. Job limits for a job on remote nodes use the `systemd` defaults and are set using the `systemd(1M)` command on the initiating host.

About Array Configuration

The system administrator has to initialize the Array configuration database, a file that is used by the Array Services daemon in executing almost every `ainfo` and `array` command. For details about array configuration, see the man pages cited in Table 1-5.

Table 1-5 Information Sources: Array Configuration

Topic	Man Page
Array Services overview	array_services(5)
Array Services user commands	ainfo(1) , array(1)
Array Services daemon overview	arrayd(1m)
Configuration file format	arrayd.conf(4) , /usr/lib/array/arrayd.conf.template
Configuration file validator	ascheck(1)
Array Services simple configurator	arrayconfig(1m)

Security Considerations for Standard Array Services

The array services daemon, `arrayd(1M)`, runs as root. As with other system services, if it is configured carelessly it is possible for arbitrary and possibly unauthorized user to disrupt or even damage a running system.

By default, most array commands are executed using the user, group, and project ID of either the user that issued the original command, or `arraysvcs`. When adding new array commands to `arrayd.conf`, or modifying existing ones, always use the most restrictive IDs possible in order to minimize trouble if a hostile or careless user were to run that command. Avoid adding commands that run with "more powerful" IDs (such as user "root" or group "sys") than the user. If such commands are necessary, analyze them carefully to ensure that an arbitrary user would not be granted any more privileges than expected, much the same as one would analyze a `setuid` program.

In the default array services configuration, the `arrayd` daemon allows all the local requests to access `arrayd` but not the remote requests. In order to let the remote requests access the `arrayd`, the `AUTHENTICATION` parameter needs to be set to `NONE` in the `/usr/lib/array/arrayd.auth` file. By default it is set to `NOREMOTE`. When the `AUTHENTICATION` parameter is set to `NONE`, the `arrayd` daemon assumes that a remote user will accurately identify itself when making a request. In other words, if a request claims to be coming from user "abc", the `arrayd` daemon assumes that it is in

fact from user "abc" and not somebody spoofing "abc". This should be adequate for systems that are behind a network firewall or otherwise protected from hostile attack, and in which all the users inside the firewall are presumed to be non-hostile. On systems, for which this is not the case (for example, those that are attached to a public network, or when individual machines otherwise cannot be trusted), the Array Services AUTHENTICATION parameter should be set to NOREMOTE. When AUTHENTICATION is set to NONE, all requests from remote systems are authenticated using a mechanism that involves private keys that are known only to the super-users on the local and remote systems. Requests originating on systems that do not have these private keys are rejected. For more details, see the section on "Authentication Information" in the `arrayd.conf(4)` man page.

The `arrayd` daemon does not support mapping user, group or project names between two different namespaces; all members of an array are assumed to share the same namespace for users, groups, and projects. Thus, if systems "A" and "B" are members of the same array, username "abc" on system A is assumed to be the same user as username "abc" on system B. This is most significant in the case of username "root". Authentication should be used if necessary to prevent access to an array by machines using a different namespace.

About the Uses of the Configuration File

The configuration files are read by the Array Services daemon when it starts. Normally it is started in each node during the system startup. (You can also run the daemon from a command line in order to check the syntax of the configuration files.)

The configuration files contain data needed by `ainfo` and `array`:

- The names of Array systems, including the current Array but also any other Arrays on which a user could run an Array Services command (reported by `ainfo`).
- The names and types of the nodes in each named Array, especially the hostnames that would be used in an Array Services command (reported by `ainfo`).
- The authentication keys, if any, that must be used with Array Services commands (required as `-Kl` and `-Kr` command options, see "Summary of Common Command Options" on page 14).
- The commands that are valid with the `array` command.

About Configuration File Format and Contents

A configuration file is a readable text file. The file contains entries of the following four types, which are detailed in later topics.

Array definition	Describes this array and other known arrays, including array names and the node names and types.
Command definition	Specifies the usage and operation of a command that can be invoked through the <code>array</code> command.
Authentication	Specifies authentication numbers that must be used to access the Array.
Local option	Options that modify the operation of the other entries or <code>arrayd</code> .

Blank lines, white space, and comment lines beginning with “#” can be used freely for readability. Entries can be in any order in any of the files read by `arrayd`.

Besides punctuation, entries are formed with a keyword-based syntax. Keyword recognition is not case-sensitive; however keywords are shown in uppercase in this text and in the man page. The entries are primarily formed from keywords, numbers, and quoted strings, as detailed in the man page `arrayd.conf(4)`.

Loading Configuration Data

The Array Services daemon, `arrayd`, can take one or more filenames as arguments. It reads them all, and treats them like logical continuations (in effect, it concatenates them). If no filenames are specified, it reads `/usr/lib/array/arrayd.conf` and `/usr/lib/array/arrayd.auth`. A different set of files, and any other `arrayd` command-line options, can be written into the file `/etc/config/arrayd.options`, which is read by the startup script that launches `arrayd` at boot time.

Since configuration data can be stored in two or more files, you can combine different strategies, for example:

- One file can have different access permissions than another. Typically, `/usr/lib/array/arrayd.conf` is world-readable and contains the available array commands, while `/usr/lib/array/arrayd.auth` is readable only by root and contains authentication codes.

- One node can have different configuration data than another. For example, certain commands might be defined only in certain nodes; or only the nodes used for interactive logins might know the names of all other nodes.
- You can use NFS-mounted configuration files. You could put a small configuration file on each machine to define the Array and authentication keys, but you could have a larger file defining `array` commands that is NFS-mounted from one node.

After you modify the configuration files, you can make `arrayd` reload them by killing the daemon and restarting it in each machine. The script `/etc/init.d/array` supports this operation:

To kill daemon, execute this command:

```
/etc/init.d/array stop
```

To kill and restart the daemon in one operation; perform the following command:

```
/etc/init.d/array restart
```

The Array Services daemon in any node knows only the information in the configuration files available in that node. This can be an advantage, in that you can limit the use of particular nodes; but it does require that you take pains to keep common information synchronized. (An automated way to do this is summarized under "Designing New Array Commands" on page 37.)

About Substitution Syntax

The man page `arrayd.conf(4)` details the syntax rules for forming entries in the configuration files. An important feature of this syntax is the use of several kinds of text substitution, by which variable text is substituted into entries when they are executed.

Most of the supported substitutions are used in command entries. These substitutions are performed dynamically, each time the `array` command invokes a subcommand. At that time, substitutions insert values that are unique to the invocation of that subcommand. For example, the value `%USER` inserts the user ID of the user who is invoking the `array` command. Such a substitution has no meaning except during execution of a command.

Substitutions in other configuration entries are performed only once, at the time the configuration file is read by `arrayd`. Only environment variable substitution makes sense in these entries. The environment variable values that are substituted are the

values inherited by `arrayd` from the script that invokes it, which is `/etc/init.d/array`.

Testing Configuration Changes

The configuration files contain many sections and options (detailed in the section that follow this one). The Array Services command `ascheck` performs a basic sanity check of all configuration files in the array.

After making a change, you can test an individual configuration file for correct syntax by executing `arrayd` as a command with the `-c` and `-f` options. For example, suppose you have just added a new command definition to `/usr/lib/array/arrayd.local`. You can check its syntax with the following command:

```
arrayd -c -f /usr/lib/array/arrayd.local
```

When testing new commands for correct operation, you need to see the warning and error messages produced by `arrayd` and processes that it may spawn. The `stderr` messages from a daemon are not normally visible. You can make them visible by the following procedure:

1. On one node, kill the daemon, as follows:

```
# /etc/init.d/array stop
```

2. In one shell window on that node, start `arrayd` with the options `-n -v`, as follows:

```
# /usr/sbin/arrayd -n -v
```

Instead of moving into the background, it remains attached to the shell terminal.

Note: Although `arrayd` becomes functional in this mode, it does not refer to `/etc/config/arrayd.options`, so you need to specify explicitly all command-line options, such as the names of nonstandard configuration files.

3. From another shell window on the same or other nodes, issue `ainfo` and `array` commands to test the new configuration data. Diagnostic output appears in the `arrayd` shell window.
4. Terminate `arrayd` and restart it as a daemon (without `-n`).

During steps 1, 2, and 4, the test node may fail to respond to `ainfo` and `array` commands, so users should be warned that the Array is in test mode.

Configuring Arrays and Machines

Each ARRAY entry gives the name and composition of an Array system that users can access. At least one ARRAY must be defined at every node, the array in use.

Note: ARRAY is a keyword.

Specifying Arrayname and Machine Names

A simple example of an ARRAY definition is as follows:

```
array simple
    machine congo
    machine niger
    machine nile
```

The arrayname `simple` is the value the user must specify in the `-a` option (see "Summary of Common Command Options" on page 14). One arrayname should be specified in a DESTINATION ARRAY local option as the default array (reported by `ainfo dflt`). Local options are listed under "Configuring Local Options" on page 36.

Note: It is recommended that you have at least one array called `me` that just contains the localhost.

The default `arrayd.conf` file has the `me` array defined as the default destination array.

The MACHINE subentries of ARRAY define the node names that the user can specify with the `-s` option. These names are also reported by the command `ainfo machines`.

Specifying IP Addresses and Ports

The simple MACHINE subentries shown in the example are based on the assumption that the hostname is the same as the machine's name to Domain Name Services (DNS). If a machine's IP address cannot be obtained from the given hostname, you must provide a HOSTNAME subentry to specify either a completely qualified domain name or an IP address, as follows:

```
array simple
  machine congo
    hostname congo.engr.hitech.com
    port 8820
  machine niger
    hostname niger.engr.hitech.com
  machine Nile
    hostname "198.206.32.85"
```

The preceding example also shows how the PORT subentry can be used to specify that arrayd in a particular machine uses a different socket number than the default 5434.

Specifying Additional Attributes

Under both ARRAY and MACHINE you can insert attributes, which are named string values. These attributes are not used by Array Services, but they are displayed by `ainfo`. Some examples of attributes would be as follows:

```
array simple
  array_attribute config_date="04/03/96"
  machine a_node
  machine_attribute aka="congo"
  hostname congo.engr.hitech.com
```

Tip: You can write code that fetches any arrayname, machine name, or attribute string from any node in the array.

Configuring Authentication Codes

In Array Services only one type of authentication is provided: a simple numeric key that can be required with any Array Services command. You can specify a single

authentication code number for each node. The user must specify the code with any command entered at that node, or addressed to that node using the `-s` option (see "Summary of Common Command Options" on page 14).

The `arshell` command is like `rsh` in that it runs a command on another machine under the `userid` of the invoking user. Use of authentication codes makes Array Services somewhat more secure than `rsh`.

Configuring Array Commands

The user can invoke arbitrary system commands on single nodes using the `arshell` command (see "Using `arshell`" on page 22). The user can also launch MPI and PVM programs that automatically distribute over multiple nodes. However, the only way to launch coordinated system programs on all nodes at once is to use the `array` command. This command does not accept any system command; it only permits execution of commands that the administrator has configured into the Array Services database.

You can define any set of commands that your users need. You have complete control over how any single Array node executes a command (the definition can be different in different nodes). A command can simply invoke a standard system command, or, since you can define a command as invoking a script, you can make a command arbitrarily complex.

Operation of Array Commands

When a user invokes the `array` command, the subcommand and its arguments are processed by the destination node specified by `-s`. Unless the `-l` option was given, that daemon also distributes the subcommand and its arguments to all other array nodes that it knows about (the destination node might be configured with only a subset of nodes). At each node, `arrayd` searches the configuration database for a `COMMAND` entry with the same name as the array subcommand.

In the following example, the subcommand `uptime` is processed by `arrayd` in node `tokyo`:

```
array -s tokyo uptime
```

When `arrayd` finds the subcommand valid, it distributes it to every node that is configured in the default array at node `tokyo`.

The **COMMAND** entry for `uptime` is distributed in this form (you can read it in the file `/usr/lib/array/arrayd.conf`).

```
command uptime          # Display uptime/load of all nodes in array
    invoke /usr/lib/array/auptime %LOCAL
```

The **INVOKE** subentry tells `arrayd` how to execute this command. In this case, it executes a shell script `/usr/lib/array/auptime`, passing it one argument, the name of the local node. This command is executed at every node, with `%LOCAL` replaced by that node's name.

Summary of Command Definition Syntax

Look at the basic set of commands distributed with Array Services (`/usr/lib/array/arrayd.conf`). Each **COMMAND** entry is defined using the subentries shown in Table 1-6. (These are described in great detail in the man page `arrayd.conf(4)`.)

Table 1-6 Subentries of a **COMMAND** Definition

Keyword	Meaning of Following Values
COMMAND	The name of the command as the user gives it to <code>array</code> .
INVOKE	A system command to be executed on every node. The argument values can be literals, or arguments given by the user, or other substitution values.
MERGE	A system command to be executed only on the distributing node, to gather the streams of output from all nodes and combine them into a single stream.
USER	The user ID under which the INVOKE and MERGE commands run. Usually given as <code>USER %USER</code> , so as to run as the user who invoked <code>array</code> .
GROUP	The group name under which the INVOKE and MERGE commands run. Usually given as <code>GROUP %GROUP</code> , so as to run in the group of the user who invoked <code>array</code> (see the <code>groups(1)</code> man page).

Keyword	Meaning of Following Values
PROJECT	The project under which the INVOKE and MERGE commands run. Usually given as PROJECT %PROJECT, so as to run in the project of the user who invoked <code>array</code> (see the <code>projects(5)</code> man page).
OPTIONS	A variety of options to modify this command; see Table 1-8.

The system commands called by INVOKE and MERGE must be specified as full pathnames, because `arrayd` has no defined execution path. As with a shell script, these system commands are often composed from a few literal values and many substitution strings. The substitutions that are supported (which are documented in detail in the `arrayd.conf(4)` man page) are summarized in Table 1-7.

Table 1-7 Substitutions Used in a COMMAND Definition

Substitution	Replacement Value
%1..%9; %ARG(<i>n</i>); %ALLARGS; %OPTARG(<i>n</i>)	Argument tokens from the user's subcommand. %OPTARG does not produce an error message if the specified argument is omitted.
%USER, %GROUP, %PROJECT	The effective user ID, effective group ID, and project of the user who invoked <code>array</code> .
%REALUSER, %REALGROUP	The real user ID and real group ID of the user who invoked <code>array</code> .
%ASH	The ASH under which the INVOKE or MERGE command is to run.
%PID(<i>ash</i>)	List of PID values for a specified ASH. %PID(%ASH) is a common use.
%ARRAY	The array name, either default or as given in the <code>-a</code> option.
%LOCAL	The hostname of the executing node.

Substitution	Replacement Value
%ORIGIN	The full domain name of the node where the <code>array</code> command ran and the output is to be viewed.
%OUTFILE	List of names of temporary files, each containing the output from one node's <code>INVOKE</code> command (valid only in the <code>MERGE</code> subentry).

The `OPTIONS` subentry permits a number of important modifications of the command execution; these are summarized in Table 1-8.

Table 1-8 Options of the `COMMAND` Definition

Keyword	Effect on Command
<code>LOCAL</code>	Do not distribute to other nodes (effectively forces the <code>-l</code> option).
<code>NEWSSESSION</code>	Execute the <code>INVOKE</code> command under a newly created <code>ASH</code> . <code>%ASH</code> in the <code>INVOKE</code> line is the new <code>ASH</code> . The <code>MERGE</code> command runs under the original <code>ASH</code> , and <code>%ASH</code> substitutes as the old <code>ASH</code> in that line.
<code>SETRUID</code>	Set both the real and effective user ID from the <code>USER</code> subentry (normally <code>USER</code> only sets the effective UID).
<code>SETRGID</code>	Set both the real and effective group ID from the <code>GROUP</code> subentry (normally <code>GROUP</code> sets only the effective GID).
<code>QUIET</code>	Discard the output of <code>INVOKE</code> , unless a <code>MERGE</code> subentry is given. If a <code>MERGE</code> subentry is given, pass <code>INVOKE</code> output to <code>MERGE</code> as usual and discard the <code>MERGE</code> output.
<code>NOWAIT</code>	Discard the output and return as soon as the processes are invoked; do not wait for completion (a <code>MERGE</code> subentry is ineffective).

Configuring Local Options

The `LOCAL` entry specifies options to `arrayd` itself. The most important options are summarized in Table 1-9.

Table 1-9 Subentries of the LOCAL Entry

Subentry	Purpose
DIR	Pathname for the <code>arrayd</code> working directory, which is the initial, current working directory of INVOKE and MERGE commands. The default is <code>/usr/lib/array</code> .
DESTINATION ARRAY	Name of the default array, used when the user omits the <code>-a</code> option. When only one ARRAY entry is given, it is the default destination.
USER, GROUP, PROJECT	Default values for COMMAND execution when USER, GROUP, or PROJECT are omitted from the COMMAND definition.
HOSTNAME	Value returned in this node by <code>%LOCAL</code> . Default is the hostname.
PORT	Socket to be used by <code>arrayd</code> .

If you do not supply LOCAL USER, GROUP, and PROJECT values, the default values for USER and GROUP are “arraysvcs.”

The HOSTNAME entry is needed whenever the `hostname` command does not return a node name as specified in the ARRAY MACHINE entry. In order to supply a LOCAL HOSTNAME entry unique to each node, each node needs an individualized copy of at least one configuration file.

Designing New Array Commands

A basic set of commands is distributed in the file `/usr/lib/array/arrayd.conf.template`. You should examine this file carefully before defining commands of your own. You can define new commands which then become available to the users of the Array system.

Typically, a new command will be defined with an INVOKE subentry that names a script written in `sh`, `csh`, or `Perl` syntax. You use the substitution values to set up arguments to the script. You use the USER, GROUP, PROJECT, and OPTIONS subentries to establish the execution conditions of the script. For one example of a command definition using a simple script, see “About the Distributed Example” on page 23.

Within the invoked script, you can write any amount of logic to verify and validate the arguments and to execute any sequence of commands. For an example of a script in Perl, see `/usr/lib/array/aps`, which is invoked by the `array ps` command.

Note: Perl is a particularly interesting choice for `array` commands, since Perl has native support for socket I/O. In principle at least, you could build a distributed application in Perl in which multiple instances are launched by `array` and coordinate and exchange data using sockets. Performance would not rival the highly tuned MPI and PVM libraries, but development would be simpler.

The administrator has need for distributed applications as well, since the configuration files are distributed over the Array. Here is an example of a distributed command to reinitialize the Array Services database on all nodes at once. The script to be executed at each node, called `/usr/lib/array/arrayd-reinit` would read as follows:

```
#!/bin/sh
# Script to reinitialize arrayd with a new configuration file
# Usage:  arrayd-reinit <hostname:new-config-file>
sleep 10      # Let old arrayd finish distributing
rcp $1 /usr/lib/array/
/etc/init.d/array restart
exit 0
```

The script uses `rcp` to copy a specified file (presumably a configuration file such as `arrayd.conf`) into `/usr/lib/array` (this will fail if `%USER` is not privileged). Then the script restarts `arrayd` (see `/etc/init.d/array`) to reread configuration files.

The command definition would be as follows:

```
command reinit
    invoke /usr/lib/array/arrayd-reinit %ORIGIN:%1
    user   %USER
    group  %GROUP
    options nowait    # Exit before restart occurs!
```

The `INVOKE` subentry calls the restart script shown above. The `NOWAIT` option prevents the daemon's waiting for the script to finish, since the script will kill the daemon.

Secure Array Services

This section provides more detailed information about Secure Array Services (SAS) and covers the following topics:

- "Differences between Standard and Secure Array Services" on page 39
- "Secure Array Services Certificates" on page 40
- "Secure Array Services Parameters" on page 44
- "Secure Shell Considerations" on page 44

Differences between Standard and Secure Array Services

This section describes the differences between standard Array Services (`arraysvcs`) and Secure Array Services (`sarraysvcs`). Secure Array Services uses Secure Sockets Layer (SSL) to communicate between `arrayd` daemons.

You need to make sure to properly protect access to various certificate files since they contain private information and keys for accessing the software.

A summary of the differences between standard Array Services and Secure Array Services is, as follows:

- You **cannot** install and run standard Array Services (`arraysvcs`) and Secure Array Services (`sarraysvcs`) on the same system.
- All the hosts in an array must run either `arraysvcs` or `sarraysvcs`. The two versions cannot operate at the same time.
- The daemon for Secure Array Services is `sarrayd`. For standard Array Services, it is `arrayd`.
- Secure Array Services requires secure shell (SSH) to be installed.
- The `chkconfig` flag is, as follows:
 - `array` for standard Array Services
 - `sarray` on Secure Array Services
- The startup script is, as follows:
 - `/etc/init.d/sarray` for Secure Array Services

- `/etc/init.d/sarray` for standard Array Services
- The `arshell(1)` command is not available for Secure Array Services.
- For Secure Array Services, all command requests are sent to the local `sarrayd` running on the current host. This is for security reasons.
- Secure Array Services requires OpenSSL images to be installed in order to run.
- On Secure Array Services, the `AUTHENTICATION` parameter is set and enforced to `NONE` since certificates are used from the server side and the client side.
- On Secure Array Services, some additional `arrayd.conf` parameters are available, as follows:

Parameter	Default setting
-----------	-----------------

<code>ssl_verify_depth</code>	
-------------------------------	--

	1
--	---

<code>ssl_cipher_list</code>	
------------------------------	--

	"ALL:!ADH:!LOW:!EXP:!MD5:@STRENGTH"
--	-------------------------------------

- The default certificate installed in `/usr/lib/array/cert` requires a local user (normally `arraysvcs`) to have group read access (normally `arraysvcs`) to the files in this directory. This means any user defined for a particular command section must have the same group access.

On Secure Array Services, is important to make sure group-read access is restricted to very few accounts. Not doing so can compromise the security features of SAS.

Secure Array Services Certificates

Certificates are used for authentication and for negotiating encryption of subsequent traffic between the `sarrayd` daemons in an array. The current implementation require the server and client certificates to be present. Upon starting, the `/etc/init.d/sarray` script attempts to create the required certificates using the `makecert` script. Certificates are not over-written. The default certificate, upon installation, allows a host to run stand-alone. The `makecert` script can be found under `/usr/lib/array/`.

Note: The first invocation of the `sarray` services may take from five to ten minutes because it has to generate the Diffie-Hellman keys required for proper certificate exchange.

If it is required to run Array Services in a cluster, you need to sign SAS root certificate common to the entire cluster (see the `gencert` command information that follows later in this section). The SAS root certificate can be self-signed (default) or signed with any valid certificate obtained from an external source.

The layout of certificate files is, as follows:

- `/usr/lib/array/cert` directory; characteristics are, as follows: permissions=750, owner=root, group=arraysvcs or the group defined in the `arrayd.conf` file. It contains the certificate and the Diffie-Hellman keys.

File	Description
<code>root.pem</code>	Array Services root certificate. Self- signed or signed by an external source (see the <code>makecert -k</code> option)
<code>client.pem</code> and <code>server.pem</code>	Client and server certificate signed by <code>root.pem</code>
<code>dh1024.pem</code>	Diffie-Hellman keys for certificate exchange

- `/usr/lib/array/cert/keys` directory; characteristics are, as follows: permissions 750, owner=root, group=arraysvcs or the group defined in `arrayd.conf`. It contains the passphrase file leading to the private keys. Note passphrase are randomly generated.
- If the group is different than `arraysvcs` or has been changed, use `makecert -X` to perform the required adjustments. See `makecert -h` for help. Note this will adjust keys and certificate ownership and permissions according to `arrayd.conf` user and group in local section. If not defined, `arraysvcs` is used for group and user. On Linux, `groupadd` and `useradd` are available from the Linux distribution and should be manually executed.
- You can generate certificate for an entire cluster by running the `gencert` command. See `gencert -h` for help.

There are two certificate utilities available. They both use `/usr/bin/openssl` command-line utilities. They reside in the `/usr/lib/array` directory and can only be executed by root. Descriptions of these certificate utilities are, as follows:

- The `makecert` utility is used to manipulate certificates for Secure Array Services. You can use the `makecert -h` command to obtain more information. A portion of the help output is, as follows:

```

makecert -C {root|serverCA|server|client|dh} ...
makecert -K
makecert -v certificate_file_name options
makecert { -i|-I } {root|serverCA|server|client|dh}
...
makecert { -g|-G } [ -P password_length ] \
                  [ -D duration_in_months ] \
                  [ -k signing_key_file ] \
                  [ -S subject_prefix ] \
                  [ -H subject_FQHN ] \
                  {root|serverCA|server|client|dh} \
                  [ signing_certificate_file ] \
                  additional_certificates_to_be_included ]
makecert -X {new_group_ownership for certificate}

where:-C   Clean directory specified on command line
-K        WARNING: perform a
          rm -rf /usr/lib/array/cert/
-i        Install/do_not_overwrite command-line
          files under /usr/lib/array/cert directory.
-I        Same as -i but overwrite files.
-g        Generate certificate or Diffie-Hellman
          keys using /dev/urandom only.
-G        Same as -g but not limited to /dev/urandom
          files.

```

- The `gencert` utility is used to generate certificates for multiple hosts, generally for an array defined in the `arrayd.auth` file. The first host on the command-line will serve as the root certificate. It uses `makecert` to generate certificates. You can use the `gencert -h` command to obtain more information. A portion of the help output is, as follows:

```

gencert [ -hdC ] [ hostnameRoot hostname ] ...
gencert [ -hdC ] -k signCertFile,signKeyFile
          hostname hostname ...
where: -h   this help message

        -d   Also generate Diffie-Hellman 1024
             bits keys
        -C   Do not clean non-local entries after

```

generation.

The `gencert` command generates all the necessary certificates for the specified hostname on the command line. For the current host, certificates are installed in `/usr/lib/array/cert`. For the remaining hosts, a tar file is created for each one.

The first utility expects the first hostname to be the base root for the remaining hosts. If no hosts are supplied, the current host is used.

The second utility is used when an external passphrase (the one that generated the signing private key) is used. In this case, `root.pem` is signed by the specified keys.

Secure Array Services Parameters

Currently, there are two parameters for Secure Array Services in the `arrayd.conf` file, as follows:

Parameter	Default
-----------	---------

<code>ssl_verify_depth</code>	
-------------------------------	--

	Default = 1. This should be changed if SAS root certificate are not self-signed.
--	--

<code>ssl_cipher_list</code>	
------------------------------	--

	Default = ALL:!ADH:!LOW:!EXP:!MD5:@STRENGTH
--	---

Secure Shell Considerations

Secure Array Services requires that some version of Secure Shell (SSH) has been installed and is running. Secure Array Services has been tested with OpenSSH versions 3.6 and 4.1, but it should work with any version of SSH. The only known requirement for the current version of Secure Array Services is that the SSH implementation installed support the following options:

Option	Description
--------	-------------

<code>-f</code>	
-----------------	--

	Requests SSH to go to background just before command execution.
--	---

-i *identity_file* Selects a file from which the identity (private key) for RSA or DSA authentication is read.

Any authentication strategy supported by SSH can be used for Secure Array Services. However, having to enter passwords for every host where application processes will execute is tedious and prone to error, as a result SGI discourages the use of such an approach for authentication. Also, executions of MPI applications launched via batch schedulers cannot be authenticated interactively. A better approach to authentication via SSH is the use of key agents, such as `ssh-agent`, or an unencrypted key. For additional information about these approaches we recommend you consult *SSH, The Secure Shell: The Definitive Guide* from O'Reilly publishing.

Since there is no standard defined location for the `ssh` client command, the full path for the desired client can be specified using the `ARRAY_SSH_PATH` environment variable. If this environment variable is undefined it is assumed that the client command is named `ssh` and resides in the defined `PATH` of the user.

SSH allows users to authenticate using multiple identities. Support for this is provided in Secure Array Services via the `ARRAY_SSH_IDENT` environment variable. If this environment variable is set the value will be used as the identity for authentication when the `ssh` client is authenticating on the remote system. The use of identities is particularly useful when different authentication methods depending upon if the user is trying to authenticate for an interactive session or a batch session. If this variable is undefined, the default identity of the user is used.

Cpusets on Linux

This chapter describes the cpuset facility on systems running SGI® Performance Suite SGI® Accelerate™ software and covers the following topics:

- "An Overview of the Advantages Gained by Using Cpusets" on page 47
- "Cpuset Permissions" on page 67
- "Cpuset File System Directories" on page 56
- "CPU Scheduling and Memory Allocation for Cpusets" on page 68
- "Using Cpusets at the Shell Prompt" on page 70
- "Cpuset Command Line Utility"
- "Boot Cpuset" on page 76
- "Configuring a User Cpuset for Interactive Sessions" on page 79
- "Cpuset Text Format" on page 82
- "Modifying the CPUs in a Cpuset and Kernel Processing" on page 83
- "Using Cpusets with Hyper-Threads" on page 84
- "Cpuset Programming Model" on page 86
- "System Error Messages" on page 88

An Overview of the Advantages Gained by Using Cpusets

The cpuset facility is primarily a workload manager tool permitting a system administrator to restrict the number of processor and memory resources that a process or set of processes may use. A *cpuset* defines a list of CPUs and memory nodes. A process contained in a cpuset may only execute on the CPUs in that cpuset and may only allocate memory on the memory nodes in that cpuset. Essentially, cpusets provide you with a CPU and memory containers or "soft partitions" within which you can run sets of related tasks. Using cpusets on an SGI Altix system improves cache locality and memory access times and can substantially improve an application's performance and runtime repeatability. Restraining all other jobs from

using any of the CPUs or memory resources assigned to a critical job minimizes interference from other jobs on the system. For example, Message Passing Interface (MPI) jobs frequently consist of a number of threads that communicate using message passing interfaces. All threads need to be executing at the same time. If a single thread loses a CPU, all threads stop making forward progress and spin at a barrier.

Cpusets can eliminate the need for a gang scheduler, provide isolation of one such job from other tasks on a system, and facilitate providing equal resources to each thread in a job. This results in both optimum and repeatable performance.

In addition to their traditional use to control the placement of jobs on the CPUs and memory nodes of a system, cpusets also provide a convenient mechanism to control the use of Hyper-Threading Technology.

Cpusets are represented in a hierarchical virtual file system. Cpusets can be nested and they have file-like permissions.

The `sched_setaffinity`, `mbind`, and `set_mempolicy` system calls allow you to specify the CPU and memory placement for individual tasks. On smaller or limited-use systems, these calls may be sufficient.

The kernel cpuset facility provides additional support for system-wide management of CPU and memory resources by related sets of tasks. It provides a hierarchical structure to the resources, with filesystem-like namespace and permissions, and support for guaranteed exclusive use of resources.

You can have a *boot* cpuset running the traditional daemon and server tasks and a second cpuset to hold interactive `telnet`, `rlogin` and/or secure shell (SSH) user sessions called the *user* cpuset.

Creating a user cpuset provides additional isolation between interactive user login sessions and essential system tasks. For example, a user process in the user cpuset consuming excessive CPU or memory resources will not seriously impact essential system services in the boot cpuset. For more information, see "Configuring a User Cpuset for Interactive Sessions" on page 79.

This section covers the following topics:

- "Linux 2.6 Kernel Support for Cpusets" on page 49
- "Cpuset Facility Capabilities" on page 50
- "Initializing Cpusets" on page 50
- "How to Determine if Cpusets are Installed" on page 51

- "Fine-grained Control within Cpusets" on page 52
- "Cpuset Interaction with Other Placement Mechanism" on page 52
- "Cpusets and Thread Placement" on page 54
- "Safe Job Migration and Cpusets" on page 54

Linux 2.6 Kernel Support for Cpusets

The Linux 2.6 kernel provides the following support for cpusets:

- Each task has a link to a cpuset structure that specifies the CPUs and memory nodes available for its use.
- Hooks in the `sched_setaffinity` system call, used for CPU placement, and in the `mbind` system call, used for memory placement, ensure that any requested CPU or memory node is available in that task's cpuset.
- All tasks sharing the same placement constraints reference the same cpuset.
- Kernel cpusets are arranged in a hierarchical virtual file system, reflecting the possible nesting of "soft partitions".
- The kernel task scheduler is constrained to only schedule a task on the CPUs in that task's cpuset.
- The kernel memory allocation mechanism is constrained to only allocate physical memory to a task from the memory nodes in that task's cpuset.
- The kernel memory allocation mechanism provides an economical, per-cpuset metric of the aggregate memory pressure of the tasks in a cpuset. *Memory pressure* is defined as the frequency of requests for a free memory page that is not easily satisfied by an available free page. For more information, see "Memory Pressure of a Cpuset" on page 63.
- The kernel memory allocation mechanism provides an option that allows you to request that memory pages used for file I/O (the kernel page cache) and associated kernel data structures for file inodes and directories be evenly spread across all the memory nodes in a cpuset. Otherwise, they are preferentially allocated on whatever memory node that the task first accessed the memory page.
- You can control the memory migration facility in the kernel using per-cpuset files. When the memory nodes allowed to a task by cpusets changes, any memory

pages no longer allowed on that node may be migrated to nodes now allowed. For more information, see "Safe Job Migration and Cpusets" on page 54.

Cpuset Facility Capabilities

A cpuset constrains the jobs (set of related tasks) running in it to a subset of the system's memory and CPUs. The cpuset facility allows you and your system service software to do the following:

- Create and delete named cpusets.
- Decide which CPUs and memory nodes are available to a cpuset.
- Attach a task to a particular cpuset.
- Identify all tasks sharing the same cpuset.
- Exclude any other cpuset from overlapping a given cpuset, thereby, giving the tasks running in that cpuset exclusive use of those CPUs and memory nodes.
- Perform bulk operations on all tasks associated with a cpuset, such as varying the resources available to that cpuset or hibernating those tasks in temporary favor of some other job.
- Perform sub-partitioning of system resources using hierarchical permissions and resource management.

Initializing Cpusets

The kernel, at system boot time, initializes one cpuset, the root cpuset, containing the entire system's CPUs and memory nodes. Subsequent user space operations can create additional cpusets.

Mounting the cpuset virtual file system (VFS) at `/dev/cpuset` exposes the kernel mechanism to user space. This VFS allows for nested resource allocations and the associated hierarchical permission model.

You can initialize and perform other cpuset operations, using any of the these three mechanisms, as follows:

- You can create, change, or query cpusets by using shell commands on `/dev/cpuset`, such as `echo(1)`, `cat(1)`, `mkdir(1)`, or `ls(1)` as described in "Using Cpusets at the Shell Prompt" on page 70.
- You can use the `cpuset(1)` command line utility to create or destroy cpusets or to retrieve information about existing cpusets and to attach processes to existing cpusets as described in "Cpuset Command Line Utility" on page 72.
- You can use the `libcpuset` C programming application programming interface (API) functions to query or change them from within your application as described in Appendix A, "Cpuset Library Functions" on page 93. You can find information about `libcpuset` at `/usr/share/doc/packages/libcpuset/libcpuset.html`.

How to Determine if Cpusets are Installed

You can issue several commands to determine whether cpusets are installed on your system, as follows:

1. Use the `grep(1)` command to search the `/proc/filesystems` for cpusets, as follows:

```
% grep cpuset /proc/filesystems
nodev cpuset
```

2. Determine if `cpuset tasks` file is present on your system by changing directory to `/dev/cpuset` and listing the content of the directory, as follows:

```
% cd /dev/cpuset
Directory: /dev/cpuset
```

```
% ls
cgroup.procs  libvirt      memory_migrate      memory_spread_page  notify_on_release  sched_relax_domain_level
cpu_exclusive mem_exclusive memory_pressure     memory_spread_slab  release_agent      tasks
cpus          mem_hardwall memory_pressure_enabled  mems                 sched_load_balance
```

3. If the `/dev/cpuset/tasks` file is not present on your system, it means the `cpuset` file system is not mounted (usually, it is automatically mounted when the system was booted). As root, you can mount the `cpuset` file system, as follows:

```
% mount -t cpuset cpuset /dev/cpuset
```

Fine-grained Control within Cpusets

Within a single cpuset, use facilities such as `taskset(1)`, `dplace(1)`, `first-touch` memory placement, `pthread`s, `sched_setaffinity` and `mbind` to manage processor and memory placement to a more fine-grained level.

The user-level bitmask library supports convenient manipulation of multiword bitmasks useful for CPUs and memory nodes. This bitmask library is required by and designed to work with the cpuset library. You can find information on the bitmask library on your system at `/usr/share/doc/packages/libbitmask/libbitmask.html`.

Cpuset Interaction with Other Placement Mechanism

The Linux 2.6 kernel supports additional processor and memory placement mechanisms, as follows:

Note: Use the `uname(1)` command to print out system information to make sure you are running the Linux 2.6.x `sn2` kernel, as follows:

```
% uname -r -s
Linux 2.6.16.14-6-default
```

- The `sched_setaffinity(2)` and `sched_getaffinity(2)` system calls set and get the CPU affinity mask of a process. This determines the set of CPUs on which the process is eligible to run. The `taskset(1)` command provides a command line utility for manipulating the CPU affinity mask of a process using these system calls. For more information, see the appropriate man page.
- The `set_mempolicy` system call sets the NUMA memory policy of the current process to *policy*. A NUMA machine has different memory controllers with different distances to specific CPUs. The memory *policy* defines in which node memory is allocated for the process.

The `get_mempolicy(2)` system retrieves the NUMA policy of the calling process or of a memory address, depending on the setting of *flags*. The `numactl(8)` command provides a command line utility for manipulating the NUMA memory policy of a process using these system calls.

- The `mbind(2)` system call sets the NUMA memory policy for the pages in a specific range of a task's virtual address space.

Cpusets are designed to interact cleanly with other placement mechanisms. For example, a batch manager can use cpusets to control the CPU and memory placement of various jobs; while within each job, these other kernel mechanisms are used to manage placement in more detail. It is possible for a batch manager to change a job's cpuset placement while preserving the internal CPU affinity and NUMA memory placement policy, without requiring any special coding or awareness by the affected job.

Most jobs initialize their placement early in their timeslot, and jobs are rarely migrated until they have been running for a while. As long as a batch manager does **not** try to migrate a job at the same time as it is adjusting its own CPU or memory placement, there is little risk of interaction between cpusets and other kernel placement mechanisms.

The CPU and memory node placement constraints imposed by cpusets always override those of these other mechanisms.

Calls to the `sched_setaffinity(2)` system call automatically mask off CPUs that are not allowed by the affected task's cpuset. If a request results in all the CPUs being masked off, the call fails with `errno` set to `EINVAL`. If some of the requested CPUs are allowed by the task's cpuset, the call proceeds as if only the allowed CPUs were requested. The unallowed CPUs are silently ignored. If a task is moved to a different cpuset, or if the CPUs of a cpuset are changed, the CPU affinity of the affected task or tasks is lost. If a batch manager needs to preserve the CPU affinity of the tasks in a job that is being moved, it should use the `sched_setaffinity(2)` and `sched_getaffinity(2)` calls to save and restore each affected task's CPU affinity across the move, relative to the cpuset. The `cpu_set_t` mask data type supported by the C library for use with the CPU affinity calls is different from the `libbitmask` bitmasks used by `libcpuset`, so some coding will be required to convert between the two, in order to calculate and preserve cpuset relative CPU affinity.

Similar to CPU affinity, calls to modify a task's NUMA memory policy silently mask off requested memory nodes outside the task's allowed cpuset, and will fail if that results in requested an empty set of memory nodes. Unlike CPU affinity, the NUMA memory policy system calls do not support one task querying or modifying another task's policy. So the kernel automatically handles preserving cpuset relative NUMA memory policy when either a task is attached to a different cpuset, or a cpuset's `mems` value setting is changed. If the old and new `mems` value sets have the same size, the cpuset relative offset of affected NUMA memory policies is preserved. If the new `mems` value is smaller, the old `mems` value relative offsets are folded onto the new `mems` value, modulo the size of the new `mems`. If the new `mems` value is larger, then just the first N nodes are used, where N is the size of the old `mems` value.

Cpsets and Thread Placement

If your job uses the placement mechanisms described in "Cpset Interaction with Other Placement Mechanism" on page 52 and operates under the control of a batch manager, you **cannot** guarantee that a migration will preserve placement done using the mechanisms. These placement mechanisms use system wide numbering of CPUs and memory nodes, not cpuset relative numbering and the job might be migrated without its knowledge while it is trying to adjust its placement. That is, between the point where an application computes the CPU or memory node on which it wants to place a thread and the point where it issues the `sched_setaffinity(2)`, `mbind(2)` or `set_mempolicy(2)` call to direct such a placement, the thread might be migrated to a different cpuset, or its cpuset changed to different CPUs or memory nodes, invalidating the CPU or memory node number it just computed.

The `libcgroup` library provides the following mechanisms to support cpuset relative thread placement that is robust even if the job is being migrated using a batch scheduler.

If your job needs to pin a thread to a single CPU, you can use the convenient `cpuset_pin` function. This is the most common case. For more information on `cpuset_pin`, see "Basic Cpset Library Functions" on page 94.

If your job needs to implement some other variation of placement, such as to specific memory nodes, or to more than one CPU, you can use the following functions to safely guard such code from placement changes caused by job migration, as follows:

- `cpuset_get_placement` (see "`cpuset_get_placement`" on page 119)
- `cpuset_equal_placement` (see "`cpuset_equal_placement`" on page 120)
- `cpuset_free_placement` (see "`cpuset_free_placement`" on page 120)

Safe Job Migration and Cpsets

Jobs that make use of cpuset aware thread pinning described in "Cpsets and Thread Placement" on page 54 can be safely migrated to a different cpuset or have the CPUs or memory nodes of the cpuset safely changed without destroying the per-thread placement done within the job.

Procedure 2-1 Safe Job Migration Between Cpsets

To safely migrate a job to a different cpuset, perform the following steps:

1. Suspend the tasks in the job by sending their process group a `SIGSTOP` signal.

2. Use the `cpuset_init_pidlist` function and related `pidlist` functions to determine the list of tasks in the job.
3. Use `sched_getaffinity(2)` to query the CPU affinity of each task in the job.
4. Create a new cpuset, under a temporary name, with the new desired CPU and memory placement.
5. Invoke `cpuset_migrate_all` function to move the job's tasks from the old cpuset to the new cpuset.
6. Use `cpuset_delete` to delete the old cpuset.
7. Use `rename(2)` on the `/dev/cpuset` based path of the new temporary cpuset to rename that cpuset to the old cpuset name.
8. Convert the results of the previous `sched_getaffinity(2)` calls to the new cpuset placement, preserving cpuset relative offset by using the `cpuset_c_rel_to_sys_cpu` and related functions.
9. Use `sched_setaffinity(2)` to reestablish the per-task CPU binding of each thread in the job.
10. Resume the tasks in the job by sending their process group a `SIGCONT` signal.

The `sched_getaffinity(2)` and `sched_setaffinity(2)` C library calls are limited by C library internals to systems with 1024 CPUs or less. To write code that will work on larger systems, you should use the `syscall(2)` indirect system call wrapper to directly invoke the underlying system call, bypassing the C library API for these calls.

The suspend and resume operation are required in order to keep tasks in the job from changing their per thread CPU placement between steps three and six. The kernel automatically migrates the per-thread memory node placement during step four. This is necessary, because there is no way for one task to modify the NUMA memory placement policy of another task. The kernel does not automatically migrate the per-thread CPU placement, as this can be handled by the user level process doing the migration.

Migrating a job from a larger cpuset (more CPUs or nodes) to a smaller cpuset will lose placement information and subsequently moving that cpuset back to a larger cpuset will **not** recover that information. This loss of CPU affinity can be avoided as described above, using `sched_getaffinity(2)` and `sched_setaffinity(2)` to save and restore the placement (affinity) across such a pair of moves. This loss of NUMA memory placement information cannot be avoided because one task (the one doing the migration) cannot save nor restore the NUMA memory placement policy of

another. So if a batch manager wants to migrate jobs without causing them to lose their `mbind(2)` or `set_mempolicy(2)` placement, it should only migrate to cpusets with at least as many memory nodes as the original cpuset.

Cpuset File System Directories

Cpusets are named, nested sets of CPUs and memory nodes. Each cpuset is represented by a directory in the cpuset virtual file system, normally mounted at `/dev/cpuset`, as described earlier.

The state of each cpuset is represented by small text files in the directory for the cpuset. These files may be read and written using traditional shell utilities such as `cat(1)` and `echo(1)` or using ordinary file access routines from programming languages, such as `open(2)`, `read(2)`, `write(2)` and `close(2)` from the C programming library.

To view the files in a cpuset that can be either read or written, perform the following commands:

```
% cd /dev/cpuset
% ls

cgroup.procs  libvirt      memory_migrate      memory_spread_page  notify_on_release  sched_relax_domain_lev
cpu_exclusive mem_exclusive memory_pressure     memory_spread_slab  release_agent      tasks
cpus         mem_hardwall memory_pressure_enabled  mems                sched_load_balance
```

Descriptions of the files in the cpuset directory are, as follows:

Cpuset Directory File

tasks

Description

List of process IDs (PIDs) of tasks in the cpuset. The list is formatted as a series of ASCII decimal numbers, each followed by a newline. A task may be added to a cpuset (removing it from the cpuset previously containing it) by writing its PID to that cpuset's tasks file (with or without a trailing newline.)

Note that only one PID may be written to the tasks file at a time. If

<code>notify_on_release</code>	<p>a string is written that contains more than one PID, all but the first are ignored.</p> <p>Flag (0 or 1) - If set (1), the <code>/sbin/cpuset_release_agent</code> binary is invoked, with the name (<code>/dev/cpuset</code> relative path) of that cpuset in <code>argv[1]</code>, when the last user of it (task or child cpuset) goes away. This supports automatic cleanup of abandoned cpusets.</p>
<code>cpus</code>	<p>List of CPUs that tasks in the cpuset are allowed to use. For a description of the format of the <code>cpus</code> file, see "List Format" on page 67. The CPUs allowed to a cpuset may be changed by writing a new list to its <code>cpus</code> file. Note, however, such a change does not take affect until the PIDs of the tasks in the cpuset are rewritten to the cpuset's <code>tasks</code> file.</p>
<code>cpu_exclusive</code>	<p>Flag (0 or 1) - The <code>cpu_exclusive</code> flag, when set, automatically defines scheduler domains. The kernel performs automatic load balancing of active threads on available CPUs more rapidly within a scheduler domain than it does across scheduler domains. By default, this flag is off (0). Newly created cpusets initially default this flag to off (0).</p>
<code>mems</code>	<p>List of memory nodes that tasks in the cpuset are allowed to use. For a description of the format of the <code>mems</code> file, see "List Format" on page 67.</p>

<code>mem_exclusive</code>	Flag (0 or 1) - The <code>mem_exclusive</code> flag, when set, automatically defines constraints for kernel internal memory allocations. Allocations of user space memory pages are strictly confined by the allocating task's cpuset. Allocations of kernel internal pages are only confined by the nearest enclosing cpuset that is marked <code>mem_exclusive</code> . By default, this flag is off (0). Newly created cpusets also initially default this flag to off (0).
<code>memory_migrate</code>	Flag (0 or 1). If set (1), memory migration is enabled. For more information, see "Memory Migration" on page 66.
<code>memory_pressure</code>	A measure of how much memory pressure the tasks in this cpuset are causing. Always has value zero (0) unless <code>memory_pressure_enabled</code> is enabled in the top cpuset. This is a read-only file. The <code>memory_pressure</code> mechanism makes it easy to detect when the job in a cpuset is running short of memory and needing to page memory out to swap. For more information, see "Memory Pressure of a Cpuset" on page 63.
<code>memory_pressure_enabled</code>	Flag (0 or 1). This file is only present in the root cpuset, normally at <code>/dev/cpuset</code> . If set (1), <code>memory_pressure</code> calculations are enabled for all cpusets in the system. For more information, see "Memory Pressure of a Cpuset" on page 63.

<code>memory_spread_page</code>	Flag (0 or 1). If set (1), the kernel page cache (file system buffers) are uniformly spread across the cpuset. For more information, see "Memory Spread" on page 65.
<code>memory_spread_slab</code>	Flag (0 or 1). If set (1), the kernel slab caches for file I/O (directory and inode structures) are uniformly spread across the cpuset. For more information, see "Memory Spread" on page 65.
<code>sched_load_balance</code>	Flag (0 or 1). If set (1, the default) the kernel will automatically load balance processes in that cpuset over the allowed CPUs in that cpuset. If cleared (0) the kernel will avoid load balancing processes in this cpuset, unless some other cpuset with overlapping CPUs has its <code>sched_load_balance</code> flag set. (since the Linux 2.6.24 release)
<code>sched_relax_domain_level</code>	Integer, between -1 and a small positive value. The <code>sched_relax_domain_level</code> controls the width of the range of CPUs over which the kernel scheduler performs immediate rebalancing of runnable tasks across CPUs. If <code>sched_load_balance</code> is disabled, then the setting of <code>sched_relax_domain_level</code> does not matter, as no such load balancing is done. If <code>sched_load_balance</code> is enabled, then the higher the value of the <code>sched_relax_domain_level</code> , the wider the range of CPUs over which immediate load balancing is

attempted. (since the Linux the 2.6.26 release)

For additional information, see the `cpuset(7)` man page.

A new file has been added to `/proc` file system, as follows:

`/proc/pid/cpuset`

For each task (PID), list its cpuset path, relative to the root of the cpuset file system. This is a read-only file.

There are two control fields used by the kernel scheduler and memory allocation mechanism to constrain scheduling and memory allocation to the allowed CPUs. These are two fields in the `status` file of each task, as follows:

`/proc/pid/status`

<code>Cpus_allowed</code>	A bit vector of CPUs on which this task may be scheduled
<code>Mems_allowed</code>	A bit vector of memory nodes on which this task may obtain memory

There are several reasons why a tasks `Cpus_allowed` and `Mems_allowed` values may differ from the the values in the `cpus` and `mems` file for that are allowed in its current cpuset, as follows:

- A task might use the `sched_setaffinity`, `mbind`, or `set_mempolicy` functions to restrain its placement to less than its cpuset.
- Various temporary changes to `cpus_allowed` values are done by kernel internal code.
- Attaching a task to a cpuset does not change its `mems_allowed` value until the next time that task needs kernel memory.
- Changing a cpuset's `cpus` value does not change the `Cpus_allowed` of the tasks attached to it until those tasks are reattached to that cpuset (to avoid a hook in the hotpath scheduler code in the kernel).

User space action is required to update a task's `Cpus_allowed` values after changing its cpuset. Use the the `cpuset_reattach` routine to perform this update after a changing the CPUs allowed to a cpuset.

- If the hotplug mechanism is used to remove all the CPUs, or all the memory nodes, in a cpuset, the tasks attached to that cpuset will have their `Cpus_allowed` or `Mems_allowed` values altered to the CPUs or memory nodes of the closest ancestor to that cpuset that is not empty.

The confines of a cpuset can be violated after a hotplug removal that empties a cpuset, until, and unless, the system's cpuset configuration is updated to accurately reflect the new hardware configuration, and in particular, to not define a cpuset that has no CPUs still online, or no memory nodes still online. The kernel prefers misplacing a task, over starving a task of essential compute resources.

There is one other condition under which the confines of a cpuset may be violated. A few kernel critical internal memory allocation requests, marked `GFP_ATOMIC`, must be satisfied immediately. The kernel may drop some request or malfunction if one of these allocations fail. If such a request cannot be satisfied within the current task's cpuset, the kernel relaxes the cpuset, and looks for memory anywhere it can find it. It is better to violate the cpuset than stress the kernel operation.

New cpubsets are created using the `mkdir` command at the shell (see "Using Cpubsets at the Shell Prompt" on page 70) or via the C programming language (see Appendix A, "Cpubset Library Functions" on page 93). Old cpubsets are removed using the `rmdir(1)` command. The above files are accessed using `read(2)` and `write(2)` system calls, or shell commands such as `cat(1)` and `echo(1)`.

The CPUs and memory nodes in a given cpuset are always a subset of its parent. The root cpuset has all possible CPUs and memory nodes in the system. A cpuset may be exclusive (CPU or memory) only if its parent is similarly exclusive.

Each task has a pointer to a cpuset. Multiple tasks may reference the same cpuset. Requests by a task, using the `sched_setaffinity(2)` system call to include CPUs in its CPU affinity mask, and using the `mbind(2)` and `set_mempolicy(2)` system calls to include memory nodes in its memory policy, are both filtered through that task's cpuset, filtering out any CPUs or memory nodes not in that cpuset. The scheduler will not schedule a task on a CPU that is not allowed in its `cpus_allowed` vector and the kernel page allocator will not allocate a page on a node that is not allowed in the requesting task's `mems_allowed` vector.

Exclusive Cpusets

If a cpuset is marked `cpu_exclusive` or `mem_exclusive`, no other cpuset, other than a direct ancestor or descendant, may share any of the same CPUs or memory nodes.

A cpuset that is `cpu_exclusive` has a scheduler (`sched`) domain associated with it. The `sched` domain consists of all CPUs in the current cpuset that are not part of any exclusive child cpusets. This ensures that the scheduler load balancing code only balances against the CPUs that are in the `sched` domain as described in "Cpuset File System Directories" on page 56 and not all of the CPUs in the system. This removes any overhead due to load balancing code trying to pull tasks outside of the `cpu_exclusive` cpuset only to be prevented by the `Cpus_allowed` mask of the task.

A cpuset that is `mem_exclusive` restricts kernel allocations for page, buffer, and other data commonly shared by the kernel across multiple users. All cpusets, whether `mem_exclusive` or not, restrict allocations of memory for user space. This enables configuring a system so that several independent jobs can share common kernel data, such as file system pages, while isolating the user allocation of each job to its own cpuset. To do this, construct a large `mem_exclusive` cpuset to hold all the jobs, and construct child, non-`mem_exclusive` cpusets for each individual job. Only a small amount of typical kernel memory, such as requests from interrupt handlers, is allowed to be taken outside even a `mem_exclusive` cpuset.

Notify on Release Flag

If the `notify_on_release` flag is enabled (1) in a cpuset, whenever the last task in the cpuset leaves (exits or attaches to some other cpuset) and the last child cpuset of that cpuset is removed, the kernel runs the `/sbin/cpuset_release_agent` command, supplying the pathname (relative to the mount point of the cpuset file system) of the abandoned cpuset. This enables automatic removal of abandoned cpusets.

The default value of `notify_on_release` in the root cpuset at system boot is disabled (0). The default value of other cpusets at creation is the current value of their parents `notify_on_release` setting.

The `/sbin/cpuset_release_agent` command is invoked, with the name (`/dev/cpuset` relative path) of that cpuset in `argv[1]` argument. This supports automatic cleanup of abandoned cpusets.

The usual contents of the `/sbin/cpuset_release_agent` command is a simple shell script, as follows:

```
#!/bin/sh
rmdir /dev/cpuset/$1
```

By default, the `notify_on_release` flag is off (0). Newly created cpusets inherit their `notify_on_release` flag setting from their parent cpuset. As with other flag values, this flag can be changed by writing an ASCII number 0 or 1 (with optional trailing newline) into the file, to clear or set the flag, respectively.

Memory Pressure of a Cpuset

The `memory_pressure` of a cpuset provides a simple per-cpuset metric of the rate that the tasks in a cpuset are attempting to free up in use memory on the nodes of the cpuset to satisfy additional memory requests. This enables batch managers, monitoring jobs running in dedicated cpusets, to efficiently detect what level of memory pressure that job is causing.

This is useful in the following situations:

- Tightly managed systems running a wide mix of submitted jobs that may choose to terminate or re-prioritize jobs trying to use more memory than allowed on the nodes to which they are assigned.
- Tightly coupled, long running, massively parallel, scientific computing jobs that will dramatically fail to meet required performance goals if they start to use more memory than allowed.

This mechanism provides a very economical way for the batch manager to monitor a cpuset for signs of memory pressure. It is up to the batch manager or other user code to decide when to take action to alleviate memory pressures.

If the `memory_pressure_enabled` flag in the top cpuset is (0), that is, it is **not** set, the kernel does not compute this filter and the per-cpuset files `memory_pressure` contain the value zero (0).

If the `memory_pressure_enabled` flag in the top cpuset is set (1), the kernel computes this filter for each cpuset in the system, and the `memory_pressure` file for each cpuset reflects the recent rate of such low memory page allocation attempts by tasks in said cpuset.

Reading the `memory_pressure` file of a cpuset is very efficient. This mechanism allows batch schedulers to poll these files and detect jobs that are causing memory stress. They can then take action to avoid impacting the rest of the system with a job that is trying to aggressively exceed its allowed memory.

Note: Unless enabled by setting `memory_pressure_enabled` in the top cpuset, `memory_pressure` is not computed for any cpuset and always reads a value of zero.

A running average per cpuset has the following advantages:

- The system load imposed by a batch scheduler monitoring this metric is sharply reduced on large systems because this meter is per-cpuset, rather than per-task or memory region and this avoids a scan of the system-wide tasklist on each set of queries.
- A batch scheduler can detect memory pressure with a single read, instead of having to read and accumulate results for a period of time because this meter is a running average, rather than an accumulating counter.
- A batch scheduler can obtain the key information, memory pressure in a cpuset, with a single read, rather than having to query and accumulate results over all the (dynamically changing) set of tasks in the cpuset because this meter is per-cpuset rather than per-task or memory region.

A simple, per-cpuset digital filter is kept within the kernel and updated by any task attached to that cpuset if it enters the synchronous (direct) page reclaim code.

The per-cpuset `memory_pressure` file provides an integer number representing the recent (half-life of 10 seconds) rate of direct page reclaims caused by the tasks in the cpuset in units of reclaims attempted per second, times 1000.

The kernel computes this value using a single-pole, low-pass recursive digital filter coded with 32-bit integer arithmetic. The value decays at an exponential rate.

Given the simple 32-bit integer arithmetic used in the kernel to compute this value, this meter works best for reporting page reclaim rates between one per millisecond (msec) and one per 32 (approximate) seconds. At constant rates faster than one per msec, it reaches maximum at values just under 1,000,000. At constant rates between one per msec and one per second, it stabilizes to a value $N*1000$, where N is the rate of events per second. At constant rates between one per second and one per 32 seconds, it is choppy, moving up on the seconds that have an event, and then decaying until the next event. At rates slower than about one in 32 seconds, it decays all the way back to zero between each event.

Memory Spread

There are two Boolean flag files per cpuset that control where the kernel allocates pages for the file system buffers and related in kernel data structures. They are called `memory_spread_page` and `memory_spread_slab`.

If the per-cpuset, `memory_spread_page` flag is set, the kernel spreads the file system buffers (page cache) evenly over all the nodes that the faulting task is allowed to use, instead of preferring to put those pages on the node where the task is running.

If the per-cpuset, `memory_spread_slab` flag is set, the kernel spreads some file system related slab caches, such as for inodes and directory entries, evenly over all the nodes that the faulting task is allowed to use, instead of preferring to put those pages on the node where the task is running.

The setting of these flags does not affect the anonymous data segment or stack segment pages of a task.

By default, both kinds of memory spreading are off, and memory pages are allocated on the node local to where the task is running, except perhaps as modified by the tasks NUMA memory policy or cpuset configuration. This is true as long as sufficient free memory pages are available.

When new cpusets are created, they inherit the memory spread settings of their parent.

Setting memory spreading causes allocations for the affected page or slab caches to ignore the task's NUMA memory policy and be spread instead. Tasks using `mbind()` or `set_mempolicy()` calls to set NUMA memory policies will not notice any change in these calls, as a result of their containing tasks memory spread settings. If memory spreading is turned off, the currently specified NUMA memory policy once again applies to memory page allocations.

Both `memory_spread_page` and `memory_spread_slab` are Boolean flag files. By default, they contain 0. This means the feature is off for the cpuset. If a 1 is written to this file, the named feature is turned on for the cpuset.

This memory placement policy is also known (in other contexts) as round-robin or interleave.

This policy can provide substantial improvements for jobs that need to place thread local data on the corresponding node, but that need to access large file system data sets that need to be spread across the several nodes in the job's cpuset in order to fit. Without this policy, especially for jobs that might have one thread reading in the data

set, the memory allocation across the nodes in the jobs cpuset can become very uneven.

Memory Migration

Normally, under the default setting of `memory_migrate`, once a page is allocated (given a physical page of main memory), that page stays on whatever node it was allocated, as long as it remains allocated, even if the cpuset's memory placement `mems` policy subsequently changes. The default setting has the `memory_migrate` flag disabled.

When memory migration is enabled in a cpuset, if the `mems` setting of the cpuset is changed, any memory page in use by any task in the cpuset that is on a memory node no longer allowed is migrated to a memory node that is allowed.

Also, if a task is moved into a cpuset with `memory_migrate` enabled, any memory pages it uses that were on memory nodes allowed in its previous cpuset, but which are not allowed in its new cpuset, are migrated to a memory node allowed in the new cpuset.

The relative placement of a migrated page within the cpuset is preserved during these migration operations if possible. For example, if the page was on the second valid node of the prior cpuset then the page will be placed on the second valid node of the new cpuset, if possible.

In order to maintain the cpuset relative position of pages, even pages on memory nodes allowed in both the old and new cpusets may be migrated. For example, if `memory_migrate` is enabled in a cpuset, and that cpuset's `mems` file is written, changing it from say memory nodes "4-7", to memory nodes "5-8", the following page migrations is done, in order, for all pages in the address space of tasks in that cpuset:

- First, migrate pages on node 7 to node 8.
- Second, migrate pages on node 6 to node 7.
- Third, migrate pages on node 5 to node 6.
- Fourth, migrate pages on node 4 to node 5.

In this example, pages on any memory node other than "4 through 7" will not be migrated. The order in which nodes are handled in a migration is intentionally chosen so as to avoid migrating memory to a node until any migrations from that node have first been accomplished.

Mask Format

The mask format is used to represent CPU and memory node bitmasks in the `/proc/pid/status` file. It is hexadecimal, using ASCII characters "0" - "9" and "a" - "f". This format displays each 32-bit word in hex (zero filled), and for masks longer than one word, uses a comma separator between words. Words are displayed in big-endian order (most significant first). And hexadecimal digits within a word are also in big-endian order. The number of 32-bit words displayed is the minimum number needed to display all bits of the bitmask, based on the size of the bitmask. An example of the mask format is, as follows:

```
00000001          # just bit 0 set
80000000,00000000,00000000  # just bit 95 set
00000001,00000000,00000000  # just bit 64 set
000000ff,00000000          # bits 32-39 set
00000000,000E3862          # bits 1,5,6,11-13,17-19 set
```

A mask with bits 0, 1, 2, 4, 8, 16, 32 and 64 set displays as `00000001,00000001,00010117`. The first "1" is for bit 64, the second for bit 32, the third for bit 16, the fourth for bit 8, the fifth for bit 4, and the "7" is for bits 2, 1 and 0.

List Format

The list format is used to represent CPU and memory node bitmasks (sets of CPU and memory node numbers) in the `/dev/cpuset` file system. It is a comma separated list of CPU or memory node numbers and ranges of numbers, in ASCII decimal. An example of list format is, as follows:

```
0-4,9          # bits 0, 1, 2, 3, 4, and 9 set
0-3,7,12-15    # bits 0, 1, 2, 3, 7, 12, 13, 14, and 15 set
```

Cpuset Permissions

The permissions of a cpuset are determined by the permissions of the special files and directories in the cpuset file system, normally mounted at `/dev/cpuset`.

For example, a task can put itself in some other cpuset (than its current one) if it can write the `tasks` file (see "Cpuset File System Directories" on page 56) for that cpuset (requires execute permission on the encompassing directories and write permission on that `tasks` file).

An additional constraint is applied to requests to place some other task in a cpuset. One task may not attach another task to a cpuset unless it has permission to send that task a signal.

A task may create a child cpuset if it can access and write the parent cpuset directory. It can modify the CPUs or memory nodes in a cpuset if it can access that cpuset's directory (execute permissions on the encompassing directories) and write the corresponding `cpus` or `mems` file (see "Cpuset File System Directories" on page 56).

It should be noted, however, that changes to the CPUs of a cpuset do not apply to any task in that cpuset until the task is reattached to that cpuset. If a task can write the `cpus` file, it should also be able to write the `tasks` file and might be expected to have permission to reattach the tasks therein (equivalent to permission to send them a signal).

There is one minor difference between the manner in which cpuset path permissions are evaluated by `libcpuset` and the manner in which file system operation permissions are evaluated by direct system calls. System calls that operate on file pathnames, such as the `open(2)` system call, rely on direct kernel support for a task's current directory. Therefore, such calls can successfully operate on files in or below a task's current directory, even if the task lacks search permission on some ancestor directory. Calls in `libcpuset` that operate on cpuset pathnames, such as the `cpuset_query()` call, rely on `libcpuset` internal conversion of all cpuset pathnames to full, root-based paths. They cannot successfully operate on a cpuset unless the task has search permission on all ancestor directories, starting with the usual cpuset mount point (`/dev/cpuset`).

CPU Scheduling and Memory Allocation for Cpusets

This section describes CPU scheduling and memory allocation for cpusets and covers these topics:

- "Linux Kernel CPU and Memory Placement Settings" on page 68
- "Manipulating Cpusets" on page 70

Linux Kernel CPU and Memory Placement Settings

The Linux kernel exposes to user space three important attributes of each task that the kernel uses to control that tasks processor and memory placement, as follows:

- The cpuset path of each task, relative to the root of the cpuset file system, is available in the file `/proc/pid/cpuset`. For each task (PID), the file lists its cpuset path relative to the root of the cpuset file system.
- The actual CPU bitmask used by the kernel scheduler to determine on which CPUs a task may be scheduled is displayed in the `Cpus_allowed` field of the file `/proc/pid/status` for that task *pid*.
- The actual memory node bitmask used by the kernel memory allocator to determine on which memory nodes a task may obtain memory is displayed in the `Mems_allowed` field of the file of the file `/proc/pid/status` for that task *pid*.

Each of the above files is read-only. You can ask the kernel to make changes to these settings by using the various cpuset interfaces and the `sched_setaffinity(2)`, `mbind(2)`, and `set_mempolicy(2)` system calls.

The `cpus_allowed` and `mems_allowed` status file values for a task may differ from the `cpus` and `mems` values defined in the cpuset directory for the task for the following reasons:

- A task might call the `sched_setaffinity`, `mbind`, or `set_mempolicy` system calls to restrain its placement to less than its cpuset.
- Various temporary changes to `cpus_allowed` status file values are done by kernel internal code
- Attaching a task to a cpuset does not change its `mems_allowed` status file value until the next time that task needs kernel memory.
- Changing the CPUs in a cpuset does not change the `cpus_allowed` status file value of the tasks attached to the cpuset until those tasks are reattached to it (to avoid a hook in the hotpath scheduler code in the kernel).

Use the `cpuset_reattach` routine to perform this update after a changing the CPUs allowed to a cpuset.

- If hotplug is used to remove all the CPUs or all the memory nodes in a cpuset, the tasks attached to that cpuset will have their `cpus_allowed` status file values or `mems_allowed` status file values altered to the CPUs or memory nodes when the closest ancestor to that cpuset is not empty.

Manipulating Cpusets

New cpusets are created using the `mkdir(1)` command (at the shell (see Procedure 2-2 on page 70) or in C programs (see Appendix A, "Cpuset Library Functions" on page 93)). Old cpusets are removed using the `rmdir(1)` commands. The `Cpus_allowed` and `Mems_allowed` status file files are accessed using `read(2)` and `write(2)` system calls or shell commands such as `cat` and `echo`.

The CPUs and memory nodes in a given cpuset are always a subset of its parent. The root cpuset has all possible CPUs and memory nodes in the system. A cpuset may be exclusive (CPU or memory) only if its parent is similarly exclusive.

Using Cpusets at the Shell Prompt

This section describes the use of cpusets using shell commands. For information on the `cpuset(1)` command line utility, see "Cpuset Command Line Utility" on page 72. For information on using the cpuset library functions, see Appendix A, "Cpuset Library Functions" on page 93.

When modifying the CPUs in a cpuset from the shell prompt, you must write the process ID (PID) of each task attached to that cpuset back into the cpuset's `tasks` file. When using the `libcpcuset` API, use the `cpuset_reattach()` routine to perform this step. The reasons for performing this step are described in "Modifying the CPUs in a Cpuset and Kernel Processing" on page 83.

Procedure 2-2 Starting a New Job within a Cpuset

In this procedure, you will create a new cpuset called `green`, assign CPUs 2 and 3 and memory node 1 to the new cpuset, and start a subshell running in the cpuset.

To start a new job and contain it within a cpuset, perform the following steps:

1. The cpuset system is created and initialized by the kernel at system boot. You allow user space access to the cpuset system by mounting the cpuset virtual file system (VFS) at `/dev/cpuset`, as follows:

```
% mkdir /dev/cpuset
% mount -t cpuset cpuset /dev/cpuset
```

Note: If the `mkdir(1)` and/or the `mount(8)` command fail, it is because they have already been performed.

2. Create the new cpuset called `green` within the `/dev/cpuset` virtual file system using the `mkdir` command, as follows:

```
% cd /dev/cpuset
```

```
% mkdir green
```

```
% cd green
```

3. Use the `echo` command to assign CPUs 2 and 3 and memory node 1 to the `green` cpuset, as follows:

```
% /bin/echo 2-3 > cpus
```

```
% /bin/echo 1 > mems
```

4. Start a task that will be the “parent process” of the new job and attach the task to the new cpuset by writing its PID to the `/dev/cpuset/tasks` file for that cpuset.

```
/bin/echo $$ > tasks
```

```
sh
```

5. The subshell `sh` is now running in the `green` cpuset.

The file `/proc/self/cpuset` shows your current cpuset, as follows:

```
% cat /proc/self/cpuset
```

```
/green
```

6. From this shell, you can `fork`, `exec` or `clone(2)` the job tasks. By default, any child task of this shell will also be in cpuset `green`. You can list the PIDs of the tasks currently in cpuset `green` by performing the following:

```
% cat /dev/cpuset/green/tasks
```

```
4965
```

```
5043
```

In this example, PID 4965 is your shell, and PID 5043 is the `cat` command itself.

Procedure 2-3 Removing a Cpuset from the `/dev/cpuset` Directory

To remove the cpuset `green` from the `/dev/cpuset` directory, perform the following:

1. Use the `rmdir` command to remove a directory from the `/dev/cpuset` directory, as follows:

```
%cd /dev/cpuset
%rmdir green
```

2. To determine if you can remove the cpuset, you can perform the `cat` command on the cpuset directory `tasks` files to ensure no PIDs are listed or within an application using `libcputset 'C'` API. You can also perform an `ls` command on the cpuset directory to ensure it has no subdirectories.

The `green` cpuset must be empty in order for you to remove it, if not a message similar to the following appears:

```
%rmdir green
rmdir: `green': Device or resource busy
```

Cpuset Command Line Utility

The `cpuset(1)` command is used to create and destroy cpusets, to retrieve information about existing cpusets, and to attach processes to cpusets. The `cpuset(1)` command line utility is not essential to the use of cpusets. This utility provides an alternative that may be convenient for some uses. Users of earlier versions of cpusets may find this utility familiar, though the details of the options have changed in order to reflect the current implementation of cpusets.

A cpuset is defined by a cpuset configuration file and a name. For a definition of the cpuset configuration file format, see "Cpuset Text Format" on page 82. The cpuset configuration file is used to list the CPUs and memory nodes that are members of the cpuset. It also contains any additional parameters required to define the cpuset. For more information on the cpuset configuration file, see "bootcpuset.conf File" on page 78.

This command automatically handles reattaching tasks to their cpuset whenever necessary, as described in the `cpuset_reattach` routine in Appendix A, "Cpuset Library Functions" on page 93.

The `cpuset` command accepts the following options:

Action Options (choose exactly one):

<code>-c <i>csname</i>, --create=<i>csname</i></code>	Creates cpuset named <i>csname</i> using the cpuset text format (see "Cpuset Text Format" on page 82)
---	---

	representation read from the commands input stream.
<code>-m <i>csname</i>, --modify=<i>csname</i></code>	Modifies the existing cpuset <i>csname</i> to have the properties in the cpuset text format (see "Cpuset Text Format" on page 82) representation read from the commands input stream.
<code>-x <i>csname</i>, --remove=<i>csname</i></code>	Removes the cpuset named <i>csname</i> . A cpuset may only be removed if there are no processes currently attached to it and the cpuset has no descendant cpusets.
<code>-d <i>csname</i>, --dump=<i>csname</i></code>	Writes a cpuset text format representation (see "Cpuset Text Format" on page 82) of the cpuset named <i>csname</i> to the commands output stream.
<code>-p <i>csname</i>, --procs=<i>csname</i></code>	Lists to the commands output stream the processes (by <code>pid</code>) attached to the cpuset named <i>csname</i> . If the <code>-r</code> option is also specified, lists the <code>pid</code> of each process attached to any descendant of cpuset <i>csname</i> .
<code>-a <i>csname</i>, --attach=<i>csname</i></code>	Attaches to the cpuset named <i>csname</i> the processes whose <code>pids</code> are read from the commands input stream, one <code>pid</code> per line.
<code>-i <i>csname</i>, --invoke=<i>csname</i></code>	Invokes a command in the cpuset named <i>csname</i> . If <code>-I</code> option is set, use that command and arguments, otherwise if the environment variable <code>\$SHELL</code> is set, use that command, otherwise, use <code>/bin/sh</code> .
<code>-w <i>pid</i>, --which=<i>pid</i></code>	Lists the name of the cpuset to which process <code>pid</code> is attached, to

<code>-s <i>csname</i>, --show=<i>csname</i></code>	the commands output stream. If <code>pid</code> is zero (0), then the full cpuset path of the current task is displayed.
<code>-R <i>csname</i>, --reattach=<i>r</i></code>	Prints to the commands output stream the names of the cpusets below cpuset <i>csname</i> . If the <code>-r</code> option is also specified, this recursively includes <i>csname</i> and all its descendants, otherwise it just includes the immediate child cpusets of <i>csname</i> . The cpuset names are printed one per line.
<code>-z <i>csname</i>, --size=<i>csname</i></code>	Reattaches each task in cpuset <i>csname</i> . This is required after changing the <code>cpus</code> value of a cpuset, in order to get the tasks already attached to that cpuset to rebind to the changed CPU placement.
<code>-F <i>flist</i>, --family=<i>flist</i></code>	Prints the size of (number of CPUs in) a cpuset to the commands output stream, as an ASCII decimal newline terminated string.
	Creates a family of non-overlapping child cpusets, given an <i>flist</i> of cpuset names and sizes (number of CPUs). Fails if the total sizes exceeds the size of the current cpuset. Enter cpuset names relative to the current cpuset, and their requested size, as alternating command line arguments. For example: <code>cpuset -F foo 2 bar 6 baz 4</code>

This creates three child cpusets named `foo`, `bar`, and `baz`, having 2, 6, and 4 CPUs, respectively.

This example will fail with an error message and a non-zero exit status if the current cpuset lacks at least 12 CPUs.

These cpuset names are relative to the current cpuset and will not collide with the cpuset names descendent from other cpusets. Hence two commands, running in different cpusets, can both create a child cpuset named `foo` without a problem.

Modifier Options (may be used in any combination):

`-r, --recursive`

When used with `-p` or `-s` option, applies to all descendants recursively of the named cpuset `csname`.

`-I cmd,`
`--invokecmd=cmd`

When used with the `-i` option, the command `cmd` is invoked, with any optional unused arguments. The following example invokes an interactive subshell in cpuset `foo`:

```
cpuset -i foo -I sh -- -i
```

The next example invokes a second cpuset command in cpuset `foo`, which then displays the full cpuset path of `foo`:

```
cpuset -i foo -I cpuset -- -w 0
```

Note: The double minus `--` is needed to end option parsing by the initial cpuset command.

`-f fname,`
`--file=fname`

Uses file named `fname` for command input or output stream, instead of `stdin` or `stdout`.

`--`
`move_tasks_from=csname1`
`--`
`move_tasks_to=csname2`

Move all tasks from cpuset `csname1` to cpuset `csname2`. Retries up to ten times to move all tasks, in case it is

racing against parallel attempts to fork or add tasks into cpuset *csname1*. Fails with non-zero exit status and an error message to stderr if unable to move all tasks out of *csname1*.

Help Option (overrides all other options):

-h, --help Displays command usage

Notes

The *csname* of "/" (slash) refers to the top cpuset, which encompasses all CPUs and memory nodes in the system. The *csname* of "." (dot) refers to the cpuset of the current task. If a *csname* begins with the "/" (slash) character, it is resolved relative to the top cpuset, otherwise it is resolved relative to the cpuset of the current task.

The 'command input stream' and 'command output stream' refer to the `stdin` (file descriptor 0) and `stdout` (file descriptor 1) of the command, unless the `-f` option is specified, in which case they refer to the file specified to `-f` option. Specifying the filename `-` to the `-f` option, as in `-f -`, is equivalent to not specifying the `-f` option at all.

Exactly **one** of the action options must be specified. They are, as follows:

-c, -m, -x, -d, -p, -a, -i, -w, -s, -R

The additional modifier options may be specified in any order. All modifier options are evaluated first, before the action option. If the help option is present, no action option is evaluated. The modifier options are, as follows:

-r, -I, -f

Boot Cpuset

You can use the `bootcpuset(8)` command to create a "boot"cpuset during the system boot that you can use to restrict the default placement of almost all UNIX processes on your system. You can use the `bootcpuset` to reduce the interference of system processes with applications running on dedicated cpusets.

The default cpuset for the `init` process, classic UNIX daemons, and user login shells is the root cpuset that contains the entire system. For systems dedicated to running particular applications, it is better to restrict `init`, the kernel daemons, and login shells to a particular set of CPUs and memory nodes called the `bootcpuset`.

This section covers the following topics:

- "Creating a Bootcpuset" on page 77
- "bootcpuset.conf File" on page 78

Creating a Bootcpuset

This section describes how to create a bootcpuset.

Procedure 2-4 Creating a Bootcpuset

To create a bootcpuset, perform the following steps:

1. Create `/etc/bootcpuset.conf` file with values to restrict system processes to the CPUs and memory nodes appropriate for your system, similar to the following:

```
cpus 0-1
mems 0
```

2. In the `/boot/efi/efi/SuSE/elilo.conf` file (or a similar path to the `elilo.conf` file), add the following string using the instructions that follow to the `append` argument for the kernel you are booting:

```
append="init=/sbin/bootcpuset"
```

For RHEL, you need to add the following string to the kernel line of the `grub.conf` file:

```
init=/sbin/bootcpuset
```

You should not directly edit the `elilo.conf` file because YaST and the install kernel tools may overwrite your changes when kernels are updated. Instead, edit the `/etc/elilo.conf` file and run the `elilo` command. This will place an updated `elilo.conf` in `/boot/efi/efi/SuSE` and the system will know about the change for new kernels or YaST runs.

3. Reboot your system.

Subsequent system reboots will restrict most processes to the bootcpuset defined in `/etc/bootcpuset.conf`.

bootcpuset.conf File

The `/etc/bootcpuset.conf` file describes what CPUs and memory nodes are to be in the bootcpuset. The kernel boot command line option `init` is used to invoke the `/sbin/bootcpuset` binary ahead of the `/sbin/init` binary, using the `elilo` syntax: `append="init=/sbin/bootcpuset"`.

When invoked with `pid==1`, the `/sbin/bootcpuset` binary does the following:

- Sets up a bootcpuset (configuration defined in the `/etc/bootcpuset.conf` file).
- Attaches itself to this bootcpuset.
- Attaches any unpinned kernel threads to it.
- Invokes an `exec` call to execute `/sbin/init`, `/etc/init`, `/bin/init` or `/bin/sh`.

A kernel thread is deemed to be unpinned (third bullet in the list above) if its `Cpus_allowed` value (as listed in that thread's `/proc/pid/status` file for the `Cpus_allowed` field) allows running on all online CPUs. Kernel threads that are restricted to some proper subset of CPUs are left untouched, under the assumption that they have a good reason to be running on those restricted CPUs. Such kernel threads as `migration` (to handle moving threads between CPUs) and `ksoftirqd` (to handle per-CPU work off interrupts) must be pinned to each CPU or each memory node.

Comments in the `/etc/bootcpuset.conf` configuration file begin with the `'#'` character and extend to the end of the line. After stripping comments, the `bootcpuset` command examines the first white space separated token on each line.

If the first token on the line matches `mems` or `mem` (case insensitive match) then the second token on the line is written to the `/dev/cpuset/boot/mems` file.

If the first token on the line matches `cpus` or `cpu` (case insensitive match), then the second token is written to the `/dev/cpuset/boot/cpus` file.

If the first token in its entirety matches (case insensitive match) `"verbose"`, the `bootcpuset` command prints a trace of its actions to the console. A typical such trace has 20 or 30 lines, detailing the steps taken by `/sbin/bootcpuset` and is useful in understanding its behavior and analyzing problems. The `bootcpuset` command ignores all other lines in the `/etc/bootcpuset.conf` configuration file.

Configuring a User Cpuset for Interactive Sessions

You can have a *boot* cpuset running the traditional daemon and server tasks and a second cpuset to hold interactive *telnet*, *rlogin* and/or secure shell (SSH) user sessions called the *user* cpuset.

Creating a user cpuset provides additional isolation between interactive user login sessions and essential system tasks. For example, a user process in the user cpuset consuming excessive CPU or memory resources will not seriously impact essential system services in the boot cpuset.

The following *init* script provides an example of how you can set up user cpuset. It runs as one of the last *init* scripts when the system is being booted.

If the system has a boot cpuset configured, the following script creates a second cpuset called *user* and place the *sshd* and *xinetd* daemons that create interactive login sessions in this new user cpuset. The user cpuset configuration is defined in the `/etc/usercpuset.conf` file.

To isolate the boot cpuset from the user cpuset, the set of *cpus* and *mems* values in the `/etc/bootcpuset.conf` file should not overlap the *cpus* and *mems* values in the `/etc/usercpuset.conf` file.

Procedure 2-5 Configuring a User Cpuset for Interactive Sessions

To implement a *usercpuset* for interactive sessions, perform the following:

1. Add the `/etc/usercpuset.conf` file.
2. Add the `/etc/init.d/usercpuset` script (see below).
3. Perform the following command:

```
% chkconfig --add usercpuset
```

Instructions for using this script are included in comments within the script, as follows:

```
#!/bin/sh
# /etc/init.d/usercpuset
#
# ----- RedHat -----
# chkconfig: 2345 057 99
# description: usercpuset: Put login sessions in user cpuset
# ----- SuSe -----
### BEGIN INIT INFO
# Provides: usercpuset
# Required-Start: sshd xinetd
# Required-Stop:
# Default-Start: 3 5
# Default-Stop: 0 1 2 6
# Description: Put login sessions in user cpuset
### END INIT INFO
# -----
#
# This init script creates a 'user' cpuset and places the login servers
# (sshd and xinetd) in that cpuset,
#
# This script presumes you also have a 'boot' cpuset configured,
# and does nothing if you don't.
#
# By using this init script, one can isolate essential daemon and
# server tasks from interactive user login sessions in separate
# cpusets.
#
# To use this usercpuset init script:
#
# 1) Using your editor, create a file /etc/init.d/usercpuset
#    containing this script.
# 2) Run the command the following command to insert this script
#    into the sequence of init scripts executed during system boot.
# SLES --> "insserv usercpuset"
# RHEL --> "chkconfig --add usercpuset"
# 3) Create a /etc/usercpuset.conf, in the "Cpuset Text Format"
#    described in the libcpuset(3) man page, describing what CPUs
#    ("cpus") and memory nodes ("mems") are to be used by the
```

```
# user cpuset.
# 4) Also configure and enable a boot cpuset, as documented in
# the bootcpuset(8) man page.
# 5) Beginning with the next system reboot, login sessions under
# either the SSH daemon (sshd) or xinetd (telnet, rlogin) will
# be started in the 'user' cpuset, while other daemons and
# system services, including the consolelogin, will be in the
# 'boot' cpuset.
# 6) If you did not do both steps (3) and (4) above, then this
# usercpuset script will do nothing, quietly, with no harm.

echo "### Creating user cpuset ###"

CPUSET_CMD=/usr/bin/cpuset

# Define the 'mems' and 'cpus' for user cpuset in this configuration file:
CONF=/etc/usercpuset.conf

USERCPUSET=/user

test -x $CPUSET_CMD || exit 5
test -r $CONF || exit 6

# Skip this if we didn't have a boot cpuset

test -d /dev/cpuset/boot || exit 7
if [ -f $CONF ]; then
    $CPUSET_CMD -c $USERCPUSET -f $CONF
fi

if [ -r /etc/rc.status ]; then
    # SuSE
    SSHD_PIDFILE=/var/run/sshd.init.pid
    status=$?
else
    # not SuSE
    SSHD_PIDFILE=/var/run/sshd.pid
    status=$?
fi

# sshd
```

```
$CPUSET_CMD -a $USERCPUSET < $SSHD_PIDFILE

# xinetd
echo $(pidof xinetd) | $CPUSET_CMD -a $USERCPUSET
```

Cpuset Text Format

Cpuset settings may be exported to and imported from text files using a text format representation of cpusets.

Permissions of files holding these text representations have no special significance to the implementation of cpusets. Rather, the permissions of the special cpuset files in the cpuset file system, normally mounted at `/dev/cpuset`, control reading and writing of and attaching to cpusets.

The text representation of cpusets is not essential to the use of cpusets. One can directly manipulate the special files in the cpuset file system. This text representation provides an alternative that may be convenient for some uses and a form for representing cpusets that users of earlier versions of cpusets will find familiar.

The exported cpuset text format has fewer directives than earlier Linux versions. Additional directives may be added in the future.

The cpuset text format supports one directive per line. Comments begin with the `'#'` character and extend to the end of line.

After stripping comments, the first white space separated token on each remaining line selects from the following possible directives:

<code>cpus</code>	Specifies which CPUs are in this cpuset. The second token on the line must be a comma-separated list of CPU numbers and ranges of numbers.
<code>mems</code>	Specify which memory nodes are in this cpuset. The second token on the line must be a comma-separated list of memory node numbers and ranges of numbers.
<code>cpu_exclusive</code>	The <code>cpu_exclusive</code> flag is set.
<code>mem_exclusive</code>	The <code>mem_exclusive</code> flag is set.

`notify_on_release` The `notify_on_release` flag is set

Additional unnecessary tokens on a line are quietly ignored. Lines containing only comments and white space are ignored.

The token `cpu` is allowed for `cpus` and `mem` for `mems`. Matching is case insensitive.

See the `libcputset` routines `cpuset_import` and `cpuset_export` to handle converting the internal `struct cpuset` representation of `cpusets` to (`export`) and from (`import`) this text representation.

For information on manipulating `cpuset` text files at the shell prompt or in shell scripts using the `cpuset(1)` command, see "Cpuset Command Line Utility" on page 72.

Modifying the CPUs in a Cpuset and Kernel Processing

In order to minimize the impact of `cpusets` on critical kernel code, such as the scheduler, and due to the fact that the Linux kernel does not support one task updating the memory placement of another task directly, the impact on a task of changing its `cpuset` CPU or memory node placement or of changing to which `cpuset` a task is attached, is subtle and is described in the following paragraphs.

When a `cpuset` has its memory nodes modified, for each task attached to that `cpuset`, the next time that the kernel attempts to allocate a page of memory for a particular task, the kernel notices the change in the task's `cpuset`, and updates its per-task memory placement to remain within the new `cpusets` memory placement. If the task was using memory policy `MPOL_BIND` and the nodes to which it was bound overlaps with its new `cpuset`, the task continues to use whatever subset of `MPOL_BIND` nodes that are still allowed in the new `cpuset`. If the task was using `MPOL_BIND` and now none of its `MPOL_BIND` nodes are allowed in the new `cpuset`, the task is essentially treated as if it was `MPOL_BIND` bound to the new `cpuset` (even though its NUMA placement, as queried by the `get_mempolicy()` routine, does not change). If a task is moved from one `cpuset` to another, the kernel adjusts the task's memory placement, as above, the next time that the kernel attempts to allocate a page of memory for that task.

When a `cpuset` has its CPUs modified, each task using that `cpuset` **does not change** its behavior automatically. In order to minimize the impact on the critical kernel scheduling code, tasks continue to use their prior CPU placement until they are rebound to their `cpuset` by rewriting their PID to the `tasks` file of their `cpuset`. If a task is moved from one `cpuset` to another, its CPU placement is updated in the same way as if the task's PID is rewritten to the `tasks` file of its current `cpuset`.

In summary, the memory placement of a task whose cpuset is changed is automatically updated by the kernel, on the next allocation of a page for that task but the processor placement is not updated until that task's PID is rewritten to the `tasks` file of its cpuset. The delay in rebinding a task's memory placement is necessary because the kernel does not support one task changing memory placement of another task. The added user level step in rebinding a task's CPU placement is necessary to avoid impacting the scheduler code in the kernel with a check for changes in a task's processor placement.

Using Cpusets with Hyper-Threads

Threading in a software application splits instructions into multiple streams so that multiple processors can act on them.

Hyper-Threading (HT) Technology, developed by Intel Corporation, provides thread-level parallelism on each processor, resulting in more efficient use of processor resources, higher processing throughput, and improved performance. One physical CPU can appear as two logical CPUs by having additional registers to overlap two instruction streams or a single processor can have dual-cores executing instructions in parallel.

In addition to their traditional use to control the placement of jobs on the CPUs and memory nodes of a system, cpusets also provide a convenient mechanism to control the use of Hyper-Threading Technology.

Some jobs achieve better performance by using both of the Hyper-Thread sides, A and B, of a processor core, and some run better by using just one of the sides, allowing the other side to idle.

Since each logical (Hyper-Threaded) processor in a core has a distinct CPU number, you can specify a cpuset that contains both sides of a processor core or a cpuset that contains just one side from a processor core.

Cpusets can be configured to include any combination of the logical CPUs in a system.

For example, the following cpuset configuration file called `cpuset.cfg` includes the A sides of an HT enabled system, along with all the memory, on the first 32 nodes (assuming 2 cores per node). The colon ':' prefixes the stride. The stride of '2' in this example means use every other logical CPU.

```
cpus 0-127:2    # the even numbered CPUs 0, 2, 4, ... 126
mems 0-63      # all memory nodes 0, 1, 2, ... 63
```

To create a cpuset called `foo` and run a job called `bar` in that cpuset, defined by the cpuset configuration file `cpuset.cfg` shown above, use the following commands:

```
cpuset -c /foo < cpuset.cfg
cpuset -i /foo -I bar
```

To specify both sides of the first 64 cores, use the following entry in your cpuset configuration file:

```
cpus 0-127
```

To specify just the B sides of the processor cores of an HT enabled system, use the following entry in your cpuset configuration file:

```
cpus 1-127:2
```

The examples above assume that the CPUs are uniformly numbered with the even numbers for the A side and odd numbers for the B side of the processors cores. This is usually the case, but not guaranteed. You can still place a job on a system that is not uniformly numbered. Currently, it involves a longer argument list to the `cpus` option, that is, you must explicitly list the desired CPUs.

If you are using a `bootcpuset` to keep other tasks confined, you do not need to create a separate cpuset with just the B side CPUs to avoid having some tasks running on the B sides of the processor cores. If there is no cpuset for the B sides of the processor cores, except the all encompassing root cpuset, and if only root can put tasks in the root cpuset, then no one other tasks can run on the B sides.

You can use the `dplace(1)` command to manage more detailed placement of job tasks within a cpuset. Since the `dplace` command numbering of CPUs is relative to the cpuset, it does not affect the `dplace` configuration. This is true in the case where the cpuset includes both sides of Hyper-Threaded cores, just one side of the Hyper-Threaded cores, or even is on a system that does not support hyperthreading.

Typically, the logical numbering of CPUs puts the even numbered CPUs on the A sides of processor cores and the odd numbered CPUs on the B sides. You can easily specify that only every other side is used using the stride suffix `":2"`, described above. If the CPU number range starts with an even number, the A sides of the processor cores are used. If the CPU range starts with an odd number, the be B sides of the processor cores are used.

Procedure 2-6 Configuring a System with Hyper-Threaded Cores

To setup a job to run only on the A sides of the system's Hyper-Threaded cores and to ensure that no other tasks run on the B sides (they remain idle), perform the following steps:

1. Define a bootcpuset to restrain the kernel, system daemon, and user login session threads to a designated set of CPUs.
2. Create a cpuset that includes on the A sides of the processors to be used for this job. (Either a system administrator or batch scheduler with root permission).
3. Make sure no cpuset is created using the B side CPUs in these processors to prevent disruptive tasks from running on the corresponding B side CPUs. (Either a system administrator or batch scheduler with root permission).

If you use a bootcpuset to confine the traditional UNIX load processes, nothing will run on the other CPUs in the system, except when those CPUs are included in a cpuset to which a job has been assigned. These CPUs are of course in the root cpuset, however, this cpuset is normally only usable by a system administrator or batch scheduler with root permissions. This prevents any user without root permission from running a task on those CPUs, unless an administrator or service with root permission allows it. For more information, see "Boot Cpuset" on page 76.

A `ps(1)` or `top(1)` invocation will show a handful of threads on unused CPUs. These are kernel threads assigned to every CPU in support of user applications running on those CPUs to handle tasks such as asynchronous file system writes and task migration between CPUs. If no application is actually using a CPU, the kernel threads on that CPU will be almost always idle.

Cpuset Programming Model

The programming model for this version of cpusets is an extension of the cpuset model provided on IRIX and earlier versions of SGI Linux.

The flat name space of earlier cpuset versions on SGI systems is extended to a hierarchical name space. This will become more important as systems become larger. The name space remains visible to all tasks on a system. Once created, a cpuset remains in existence until it is deleted or until the system is rebooted, even if no tasks are currently running in that cpuset.

The key properties of a cpuset are its pathname, the list of which CPUs and memory nodes it contains, and whether the cpuset has exclusive rights to these resources.

Every task (process) in the system is attached to (running inside) a cpuset. Tasks inherit their parents cpuset attachment when forked. This binding of task to a cpuset can subsequently be changed, either by the task itself, or externally from another task, given sufficient authority.

Tasks have their CPU and memory placement constrained to whatever their containing cpuset allows. A cpuset may have exclusive rights to its CPUs and memory, which provides certain guarantees that other cpusets will not overlap.

At system boot, a top level root cpuset is created, which includes all CPUs and memory nodes on the system. The usual mount point of the cpuset file system and therefore the usual file system path to this root cpuset, is `/dev/cpuset`.

Optionally, a "boot" cpuset may be created, at `/dev/cpuset/boot`, to include typically just a one or a few CPUs and memory nodes. A typical use for a "boot" cpuset is to contain the general purpose UNIX daemons and login sessions, while reserving the rest of the system for running specific major applications on dedicated cpusets. For more information, see "Boot Cpuset" on page 76.

Moved tasks do not have the memory they might have allocated on their old nodes moved to the new nodes. On kernels that support such memory migration, use the [optional] `cpuset_migrate` to move allocated memory as well.

Cpusets have a permission structure which determines which users have rights to query, modify, and attach to any given cpuset. Rights are based on the hierarchical model provided by the underlying Linux 2.6 kernel cpuset file system.

To create a cpuset from within a C language application, your program obtains a handle to a new `struct cpuset`, sets the desired attributes via that handle, and issues a `cpuset_create()` call to create the desired cpuset and bind it to the specified name. Your program can also issue calls to list by name what cpusets exist, query their properties, move tasks between cpusets, list what tasks are currently attached to a cpuset, and delete cpusets.

The names of cpusets in this C library are always relative to the root cpuset mount point, typically `/dev/cpuset`. For more information on the `libcpuset` C language application programming interface (API) functions, see Appendix A, "Cpuset Library Functions" on page 93.

System Error Messages

The Linux kernel implementation of cpusets sets `errno` to specify the reason for a failed system call that affects cpusets. These `errno` values are available when a `cpuset` library call fails. They can be displayed by shell commands used to directly manipulate files below the `/dev/cpuset` directory and can be displayed by the `cpuset(1)` command.

The possible `errno` settings and their meaning when set on a failed `cpuset` call are, as follows:

<code>ENOSYS</code>	Invoked on an operating system kernel that does not support cpusets.
<code>ENODEV</code>	Invoked on a system that supports cpusets, but when the <code>cpuset</code> file system is not currently mounted at <code>/dev/cpuset</code> .
<code>ENOMEM</code>	Insufficient memory is available.
<code>EBUSY</code>	Attempted <code>cpuset_delete()</code> on a <code>cpuset</code> with attached tasks.
<code>EBUSY</code>	Attempted <code>cpuset_delete()</code> on a <code>cpuset</code> with child cpusets.
<code>ENOENT</code>	Attempted <code>cpuset_create()</code> in a parent <code>cpuset</code> that does not exist.
<code>EEXIST</code>	Attempted <code>cpuset_create()</code> for a <code>cpuset</code> that already exists.
<code>E2BIG</code>	Attempted a <code>write(2)</code> system call on a special <code>cpuset</code> file with a length larger than some kernel determined upper limit on the length of such writes.
<code>ESRCH</code>	Attempted to <code>cpuset_move()</code> a nonexistent task.
<code>EACCES</code>	Attempted to <code>cpuset_move()</code> a task that the process lacks permission to move.
<code>ENOSPC</code>	Attempted to <code>cpuset_move()</code> a task to an empty <code>cpuset</code> .
<code>EINVAL</code>	The <code>relcpu</code> argument to <code>cpuset_pin()</code> function is out of range (not between "zero" and " <code>cpuset_size() - 1</code> ").

EINVAL	Attempted to change a cpuset in a way that would violate a <code>cpu_exclusive</code> or <code>mem_exclusive</code> attribute of that cpuset or any of its siblings.
EINVAL	Attempted to write an empty <code>cpus</code> or <code>mems</code> bitmask to the kernel. The kernel creates new cpubits (using the <code>mknod</code> function) with empty <code>cpus</code> and <code>mems</code> files and the user level cpuset and bitmask code works with empty masks. But the kernel will not allow an empty bitmask (no bits set) to be written to the special <code>cpus</code> or <code>mems</code> files of a cpuset.
EIO	Attempted to <code>write(2)</code> a string to a cpuset tasks file that does not begin with an ASCII decimal integer.
ENOSPC	Attempted to <code>write(2)</code> a list to a <code>cpus</code> file that did not include any online CPUs.
ENOSPC	Attempted to <code>write(2)</code> a list to a <code>mems</code> file that did not include any online memory nodes.
EACCES	Attempted to add a CPUS or memory resource to a cpuset that is not already in its parent.
EACCES	Attempted to set the <code>cpu_exclusive</code> or <code>mem_exclusive</code> flag on a cpuset whose parent lacks the same setting.
EBUSY	Attempted to remove a CPU or memory resource from a cpuset that is also in a child of that cpuset.
EFAULT	Attempted to read or write a cpuset file using a buffer that was outside your accessible address space.
ENAMETOOLONG	Attempted to read a <code>/proc/<i>pid</i>/cpuset</code> file for a cpuset path that was longer than the kernel page size.

NUMA Tools

The `dlook(1)` and `dplace(1)` commands are tools that can be used to improve the performance of processes running on your SGI nonuniform memory access (NUMA) machine.. You can use `dlook(1)` to find out where in memory the operating system is placing your application's pages and how much system and user CPU time it is consuming. You can use the `dplace(1)` command to bind a related set of processes to specific CPUs or nodes to prevent process migration. This can improve the performance of your application since it increases the percentage of memory accesses that are local.

The `taskset(1)` command is used to set or retrieve the CPU affinity of a running process given its PID or to launch a new command with a given CPU affinity. CPU affinity is a scheduler property that "bonds" a process to a given set of CPUs on the system. The Linux scheduler will honor the given CPU affinity and the process will not run on any other CPUs. Note that the Linux scheduler also supports natural CPU affinity; the scheduler attempts to keep processes on the same CPU as long as practical for performance reasons. Therefore, forcing a specific CPU affinity is useful only in certain applications.

Note: Information about these commands and memory locality and application performance, in general, can be found in the *Linux Application Tuning Guide*.

Cpuset Library Functions

This appendix describes the `libcpuset` C programming application programming interface (API) functions and covers the following topics:

- "Extensible Application Programming Interface" on page 93
- "Basic Cpuset Library Functions" on page 94
- "Advanced Cpuset Library Functions" on page 97

Extensible Application Programming Interface

In order to provide for the convenient and robust extensibility of this cpuset API over time, the following function enables dynamically obtaining pointers for optional functions by name, at runtime:

```
void *cpuset_function(const char * function_name)
```

It returns a function pointer or NULL if function name is not recognized.

For maximum portability, you should not reference any optional cpuset function by explicit name.

However, if you presume that an optional function will always be available on the target systems of interest, you might decide to explicitly reference it by name, in order to improve the clarity and simplicity of the software in question.

Also to support robust extensibility, flags and integer option values have names dynamically resolved at runtime, not via preprocessor macros.

Some functions in Advanced Cpuset Library Functions are marked `[optional]`. (see "Advanced Cpuset Library Functions" on page 97). They are not available in all implementations of `libcpuset`. Additional `[optional]` `cpuset_*` functions may also be added in the future. Functions that are not marked `[optional]` are available on all implementations of `libcpuset.so` and can be called directly without using `cpuset_function()`. However, any of them can also be called indirectly via `cpuset_function()`.

To safely invoke an optional function, such as for example `cpuset_migrate()`, use the following call sequence:

```
/* fp has function signature of pointer to cpuset_migrate() */
int (*fp)(struct cpuset *fromcp, struct cpuset *to_cp, pid_t pid);
fp = cpuset_function("cpuset_migrate");
if (fp) {
    fp( ... );
} else {
    puts ("cpuset migration not supported");
}
```

If you invoke an [optional] function directly, your resulting program will not be able to link with any version of `libcputset.so` that does not define that particular function.

Basic Cpuset Library Functions

The basic cpuset API provides functions usable from a C program for the processor and memory placement within a cpuset.

These functions enable an application to place various threads of its execution on specific CPUs within its current cpuset and perform related functions, such as, asking how large the current cpuset is and on which CPU within the current cpuset a thread is currently executing.

These functions do not provide the full power of the advanced cpuset API, but they are easier to use, and provide some common needs of multithreaded applications.

Unlike the rest of this document, the functions described in this section use cpuset relative numbering. In a cpuset of N CPUs, the relative cpu numbers range from zero to N - 1.

Memory placement is done automatically, preferring the node local to the requested CPU. Threads may only be placed on a single CPU. This avoids the need to allocate and free the bitmasks required to specify a set of several CPUs. These functions do not support creating or removing cpusets, only the placement of threads within an existing cpuset. This avoids the need to explicitly allocate and free cpuset structures. Operations only apply to the current thread, avoiding the need to pass the process ID of the thread to be affected.

If more powerful capabilities are needed, use the Advanced Cpuset library functions (see "Advanced Cpuset Library Functions" on page 97). These basic functions do not provide any essential new capability. They are implemented using the advanced function and are fully interoperable with them.

On error, these routines return -1 and set `errno`. If invoked on an operating system kernel that does not support cpusets, `errno` is set to `ENOSYS`. If invoked on a system that supports cpusets, but when the cpuset file system is not currently mounted at `/dev/cpuset`, the `errno` is set to `ENODEV`.

The following inclusion and linkage provides access to the cpuset API from C code:

```
#include <cpuset.h>
/* link with -lcpuset */
```

The following functions are supported in the basic cpuset C API:

<code>cpuset_pin</code>	Pins the current thread to a CPU, preferring local memory
<code>cpuset_size</code>	Returns the number of CPUs that are in the current tasks cpuset
<code>cpuset_where</code>	Returns on which CPU in current tasks cpuset did the task most recently execute
<code>cpuset_unpin</code>	Removes the affect of <code>cpuset_pin</code> , lets the task have run of its entire cpuset

`cpuset_pin`

```
int cpuset_pin(int relcpu);
```

Pins the current task to execute only on the CPU `relcpu`, which is a relative CPU number within the current cpuset of that task. Also, automatically pins the memory allowed to be used by the current task to prefer the memory on that same node (as determined by the `cpuset_cpu2node` function), but to allow any memory in the cpuset if no free memory is readily available on the same node.

Return 0 on success, -1 on error. Errors include `relcpu` being too large (greater than `cpuset_size() - 1`). They also include running on a system that does not support cpusets (`ENOSYS`) and running when the cpuset file system is not mounted at `/dev/cpuset` (`ENODEV`).

cpuset_size

```
int cpuset_size();
```

Returns the number of CPUs in the current tasks cpuset. The relative CPU numbers that are passed to the `cpuset_pin` function and that are returned by the `cpuset_where` function, must be between 0 and N - 1 inclusive, where N is the value returned by `cpuset_size`.

Returns -1 and sets `errno` on error. Errors include running on a system that does not support cpusets (ENOSYS) and running when the cpuset file system is not mounted at `/dev/cpuset` (ENODEV).

cpuset_where

```
int cpuset_where();
```

Returns the CPU number, relative to the current tasks cpuset, of the CPU on which the current task most recently executed. If a task is allowed to execute on more than one CPU, then there is no guarantee that the task is still executing on the CPU returned by `cpuset_where`, by the time that the user code obtains the return value.

Returns -1 and sets `errno` on error. Errors include running on a system that does not support cpusets (ENOSYS) and running when the cpuset file system is not mounted at `/dev/cpuset` (ENODEV).

cpuset_unpin

```
int cpuset_unpin();
```

Remove the CPU and Memory pinning affects of any previous `cpuset_pin` call, allowing the current task to execute on any CPU in its current cpuset and to allocate memory on any memory node in its current cpuset. Return -1 on error, 0 on success.

Returns -1 and sets `errno` on error. Errors include running on a system that does not support cpusets (ENOSYS) and running when the cpuset file system is not mounted at `/dev/cpuset` (ENODEV).

Advanced Cpuset Library Functions

The advanced cpuset API provides functions usable from a C language application for managing cpusets on a system-wide basis.

These functions primarily deal with the following three entities:

- `struct cpuset * structure`
- system cpusets
- tasks

The `struct cpuset * structure` provides a transient in-memory structure used to build up a description of an existing or desired cpuset. These structures can be allocated, freed, queried, and modified.

Actual kernel cpusets are created under the `/dev/cpuset` directory, which is the usual mount point of the kernel's virtual cpuset filesystem. These cpusets are visible to all tasks in the system, and persist until the system is rebooted or until the cpuset is explicitly deleted. These cpusets can be created, deleted, queried, modified, listed, and examined.

Every task (also known as a process) is bound to exactly one cpuset at a time. You can list which tasks are bound to a given cpuset, and to which cpuset a given task is bound. You can change to which cpuset a task is bound.

The primary attributes of a cpuset are its lists of CPUs and memory nodes. The scheduling affinity for each task, whether set by default or explicitly by the `sched_setaffinity()` system call, is constrained to those CPUs that are available in that task's cpuset. The NUMA memory placement for each task, whether set by default or explicitly by the `mbind()` system call, is constrained to those memory nodes that are available in that task's cpuset. This provides the essential purpose of cpusets - to constrain the CPU and Memory usage of tasks to specified subsets of the system.

The other essential attribute of a cpuset is its pathname beneath `/dev/cpuset`. All tasks bound to the same cpuset pathname can be managed as a unit, and this hierarchical name space describes the nested resource management and hierarchical permission space supported by cpusets. Also, this hierarchy is used to enforce strict exclusion, using the following rules:

- A cpuset may only be marked strictly exclusive for CPU or memory if its parent is also.

- A cpuset may not make any CPUs or memory nodes available that are not also available in its parent.
- If a cpuset is exclusive for CPU or memory, it may not overlap CPUs or memory with any of its siblings.

The combination of these rules enables checking for strict exclusion just by making various checks on the parent, siblings, and existing child cpusets of the cpuset being changed, without having to check all cpusets in the system.

On error, some of these routines return -1 or NULL and set `errno`. If one of the routines below that requires cpuset kernel support or the cpuset file system mounted is invoked on an operating system kernel that does not support cpusets, then that routine returns failure and `errno` is set to `ENOSYS`. If invoked on a system that supports cpusets, but when the cpuset file system is not currently mounted at `/dev/cpuset`, it returns failure and `errno` is set to `ENODEV`.

The following inclusion and linkage provides access to the cpuset API from C code:

```
#include <bitmask.h>
#include <cpuset.h>
/* link with -lcpuset */
```

The following functions are supported in the advanced cpuset C API:

Allocate and free struct cpuset * structure

- `cpuset_alloc` - Returns handle to newly allocated struct cpuset * structure
- `cpuset_free` - Discards no longer needed struct cpuset * structure

Lengths of CPUs and memory nodes bitmasks - needed to allocate bitmasks

- `cpuset_cpus_nbits` - Number of bits needed for a CPU bitmask on current system
- `cpuset_mems_nbits` - Number of bits needed for a memory bitmask on current system

Set various attributes of a struct cpuset * Structure

- `cpuset_setcpus` - Specifies CPUs in cpuset
- `cpuset_setmems` - Specifies memory nodes in cpuset
- `cpuset_set_iopt` - Specifies an integer value option of cpuset

- `cpuset_set_sopt` - [optional] Specifies a string value option of `cpuset`

Query various attributes of a `struct cpuset` * Structure

- `cpuset_getcpus` - Queries CPUs in `cpuset`
- `cpuset_getmems` - Queries memory nodes in `cpuset`
- `cpuset_cpus_weight` - Number of CPUs in a `cpuset`
- `cpuset_mems_weight` - Number of memory nodes in a `cpuset`
- `cpuset_get_iopt` - Query an integer value option of `cpuse`
- `cpuset_set_sopt` - [optional] Species a string value option of `cpuset`

Local CPUs and memory nodes

- `cpuset_localcpus` - Queries the CPUs local to specified memory nodes
- `cpuset_localmems` - Queries the memory nodes local to specified CPUs
- `cpuset_cpumemdist` - [optional] Hardware distance from CPU to memory node
- `cpuset_cpu2node` - Returns number of memory node closed to specified CPU
- `cpuset_addr2node` - Return number of memory node holding page at specified address.

Create, delete, query, modify, list, and examine cpusets

- `cpuset_create` - Creates a named `cpuset` as specified by `struct cpuset` * structure
- `cpuset_delete` - Deletes the specified `cpuset` (if empty)
- `cpuset_query` - Sets the `struct cpuset` structure to settings of specified `cpuset`
- `cpuset_modify` - Modifies the settings of a `cpuset` to those specified in a `struct cpuset` structure
- `cpuset_getcpusetpath` - Gets path of a tasks (0 for current) `cpuset`
- `cpuset_cpusetofpid` - Sets the `struct cpuset` structure to settings of `cpuset` of specified task
- `cpuset_mountpoint` - Returns path at which `cpuset` filesystem is mounted

- `cpuset_collides_exclusive` - [optional] True, if it would collide an exclusive

List tasks (pids) currently attached to a cpuset

- `cpuset_init_pidlist` - Initializes a list of tasks (pids) attached to a cpuset
- `cpuset_pidlist_length` - Returns number of elements in a list of pid
- `cpuset_get_pidlist` - Returns i'th element of a list of pids
- `cpuset_free_pidlist` - Deallocates a list of pids

Attach tasks to cpusets

- `cpuset_move` - Moves task (0 for current) to a cpuset
- `cpuset_move_all` - Moves all tasks in a list of pids to a cpuset
- `cpuset_migrate` - [optional] Moves a task and its memory to a cpuset
- `cpuset_migrate_all` - [optional] Moves all tasks with memory in a list of pids to a cpuset
- `cpuset_reattach` - Rebinds `cpus_allowed` of each task in a cpuset after changing its cpus

Determine memory pressure

- `cpuset_open_memory_pressure` - [optional] Opens handle to read memory_pressure
- `cpuset_read_memory_pressure` - [optional] Reads cpuset current memory_pressure
- `cpuset_close_memory_pressure` - [optional] Closes handle to read memory_pressure

Map between cpuset relative and system-wide CPU and memory node numbers

- `cpuset_c_rel_to_sys_cpu` - Maps cpuset relative CPU number to system wide number
- `cpuset_c_sys_to_rel_cpu` - Maps system-wide CPU number to cpuset relative number
- `cpuset_c_rel_to_sys_mem` - Maps cpuset relative memory node number to system wide number

- `cpuset_c_sys_to_rel_mem` - Maps system-wide memory node number to cpuset relative number
- `cpuset_p_rel_to_sys_cpu` - Maps task cpuset relative CPU number to system wide number
- `cpuset_p_sys_to_rel_cpu` - Maps system-wide CPU number to task cpuset relative number
- `cpuset_p_rel_to_sys_mem` - Maps task cpuset relative memory node number to system-wide number
- `cpuset_p_sys_to_rel_mem` - Maps system-wide memory node number to task cpuset relative number

Placement operations for detecting cpuset migration

- `cpuset_get_placement` - [optional] Returns the current placement of task `pid`
- `cpuset_equal_placement` - [optional] True, if two placements are equal
- `cpuset_free_placement` - [optional] Free placement

Bind to a CPU or memory node within the current cpuset

- `cpuset_cpupbind` - Binds to a single CPU within a cpuset (uses `sched_setaffinity(2)`)
- `cpuset_latestcpu` - Most recent CPU on which a task has executed
- `cpuset_membind` - Binds to a single memory node within a cpuset (uses `set_mempolicy(2)`)

Export cpuset settings to a regular file and import them from a regular file

- `cpuset_export` - Exports cpuset settings to a text file
- `cpuset_import` - Imports cpuset settings from a text file

Support calls to [optional] `cpuset_*` API routines

- `cpuset_function` - Returns pointer to a `libcputset.so` function, or NULL

Cpuset Library Functions Calling Sequence

A typical calling sequence would use the above functions in the following order to create a new cpuset named "xyz" and attach itself to it, as follows:

```
struct cpuset *cp = cpuset_alloc();
various cpuset_set*(cp, ...) calls
cpuset_create(cp, "xyz");
cpuset_free(cp);
cpuset_move(0, "xyz");
```

Note: Some functions are marked [optional]. For an explanation, see "Extensible Application Programming Interface" on page 93.

cpuset_alloc

```
struct cpuset *cpuset_alloc();
```

Creates, initializes, and returns a handle to a `struct cpuset` structure, that is an opaque data structure used to describe a cpuset.

After obtaining a `struct cpuset` handle with this call, you can use the various `cpuset_set()` methods to specify which CPUs and memory nodes are in the cpuset and other attributes. Then, you can create such a cpuset with the `cpuset_create()` call and free cpuset handles with the `cpuset_free()` call.

The `cpuset_alloc` function returns a zero pointer (NULL) and sets `errno` in the event that `malloc(3)` fails. See the `malloc(3)` man page for possible values of `errno` (`ENOMEM` being the most likely).

The `cpuset_alloc()` call applies a hidden undefined mark to each attribute of the allocated `struct cpuset`. Calls to the various `cpuset_set*()` routines mark the attribute being set as defined. Calls to `cpuset_create()` and `cpuset_modify()` only set the attributes of the cpuset marked defined. This is primarily noticeable when creating a cpuset. Code in the kernel sets some attributes of new cpusets, such as `memory_spread_page`, `memory_spread_slab`, and `notify_on_release`, by default to the value inherited from their parent. Unless the application using `libcpuset` explicitly overrides the setting of these attributes in the `struct cpuset`, between the calls to `cpuset_alloc()` and `cpuset_create()`, the kernel default settings will prevail. These hidden marks have no noticeable affect when modifying an existing cpuset using the sequence of calls `cpuset_alloc()`, `cpuset_query()`,

and `cpuset_modify()`, because the `cpuset_query()` call sets all attributes and marks them defined, while reading the attributes from the `cpuset`.

`cpuset_free`

```
struct cpuset *cpuset_alloc();
```

Frees the memory associated with a `struct cpuset` handle, that must have been returned by a previous `cpuset_alloc()` call. If `cp` is `NULL`, no operation is performed.

`cpuset_cpus_nbits`

```
int cpuset_cpus_nbits();
```

Return the number of bits in a CPU bitmask on current system. Useful when using `bitmask_alloc()` call to allocate a CPU mask. Some other routines below return `cpuset_cpus_nbits()` as an out-of-bounds indicator.

`cpuset_mems_nbits`

```
int cpuset_mems_nbits();
```

Returns the number of bits in a memory node bitmask on current system. Useful when using a `bitmask_alloc()` call to allocate a memory mode mask. Some other routines below return `cpuset_mems_nbits()` as an out-of-bounds indicator.

`cpuset_setcpus`

```
int cpuset_setcpus(struct cpuset *cp, const struct bitmask *cpus);
```

Given a bitmask of CPUs, the `cpuset_setcpus()` call sets the specified `cpuset cp` to include exactly those CPUs.

Returns 0 on success, else -1 on error, setting `errno`. This routine can fail if `malloc(3)` fails. See the `malloc(3)` man page for possible values of `errno` (`ENOMEM` being the most likely).

cpuset_setmems

```
void cpuset_setmems(struct cpuset *cp, const struct bitmask
*mems);
```

Given a bitmask of memory nodes, the `cpuset_setmems()` call sets the specified `cpuset cp` to include exactly those memory nodes.

Returns 0 on success, else -1 on error, setting `errno`. This routine can fail if `malloc(3)` fails. See the `malloc(3)` man page for possible values of `errno` (`ENOMEM` being the most likely).

cpuset_set_iopt

```
int cpuset_set_iopt(struct cpuset *cp, const char *optionname,
int value);
```

Sets `cpuset` integer valued option `optionname` to specified integer value. Returns 0 if `optionname` is recognized and `value` is an allowed value for that option. Returns -1 if `optionname` is recognized, but `value` is not allowed. Returns -2 if `optionname` is not recognized. Boolean options accept any non-zero value as equivalent to a value of one (1).

The following `optionname` values are recognized:

- `cpu_exclusive` - Sibling cpusets not allowed to overlap cpus (see "Exclusive Cpusets" on page 62)
- `mem_exclusive` - Sibling cpusets not allowed to overlap mems (see "Exclusive Cpusets" on page 62)
- `notify_on_release` - Invokes `/sbin/cpuset_release_agent` when `cpuset` released (see "Notify on Release Flag" on page 62)
- `memory_migrate` - Causes memory pages to migrate to new mems (see "Memory Migration" on page 66)
- `memory_spread_page` - Causes kernel buffer (page) cache to spread over `cpuset` (see "Memory Spread" on page 65)
- `memory_spread_slab` - Causes kernel file I/O data (directory and inode slab caches) to spread over `cpuset` (see "Memory Spread" on page 65)

cpuset_set_sopt

```
int cpuset_set_sopt(struct cpuset *cp, const char *optionname,  
const char *value);
```

Sets cpuset string valued option `optionname` to specified string value.

Returns 0 if `optionname` is recognized and `value` is an allowed value for that option. Returns -1 if `optionname` is recognized, but `value` is not allowed. Returns -2 if `optionname` is not recognized.

This is an [optional] function. Use the `cpuset_function()` to invoke it, as follows:

```
/* fp has function signature of pointer to cpuset_set_sopt() */  
int (*fp)(struct cpuset *cp, const char *optionname, const char *value);  
fp = cpuset_function("cpuset_set_sopt");  
if (fp) {  
    fp( ... );  
} else {  
    puts ("cpuset_set_sopt not supported");  
}
```

cpuset_getcpus

```
int cpuset_getcpus(const struct cpuset *cp, struct bitmask  
*cpus);
```

Queries CPUs in cpuset `cp`, by writing them to the bitmask `cpus`. Pass `cp == NULL` to query the current tasks cpuset.

If the memory nodes have not been set in cpuset `cp`, then no operation is performed, -1 is returned, and `errno` is set to `EINVAL`.

Returns 0 on success, else -1 on error, setting `errno`. This routine can fail if `malloc(3)` fails. See the `malloc(3)` man page for possible values of `errno` (`ENOMEM` being the most likely).

cpuset_getmems

```
int cpuset_getmems(const struct cpuset *cp, struct bitmask  
*mems);
```

Queries memory nodes in cpuset `cp`, by writing them to the bitmask `mems`. Pass `cp == NULL` to query the current tasks cpuset.

If the memory nodes have not been set in cpuset `cp`, then no operation is performed, -1 is returned, and `errno` is set to `EINVAL`.

Returns 0 on success, else -1 on error, setting `errno`. This routine can fail if `malloc(3)` fails. See the `malloc(3)` man page for possible values of `errno` (`ENOMEM` being the most likely).

`cpuset_cpus_weight`

```
int cpuset_cpus_weight(const struct cpuset *cp);
```

Queries the number of CPUs in cpuset `cp`. Pass `cp == NULL` to query the current tasks cpuset.

If the CPUs have not been set in cpuset `cp`, then zero(0) is returned.

`cpuset_mems_weight`

```
int cpuset_mems_weight(const struct cpuset *cp);
```

Queries the number of memory nodes in cpuset `cp`. Pass `cp == NULL` to query the current tasks cpuset.

If the memory nodes have not been set in cpuset `cp`, then zero (0) is returned.

`cpuset_get_iopt`

```
int cpuset_get_iopt(const struct cpuset *cp, const char *optionname);
```

Queries the value of integer option `optionname` in cpuset `cp`.

Returns value of `optionname` is recognized, else returns -1. Integer values in an uninitialized cpuset have value 0. The following `optionname` values are recognized:

- `cpu_exclusive` - Sibling cpusets not allowed to overlap cpus (see "Exclusive Cpusets" on page 62)

- `mem_exclusive` - Sibling cpusets not allowed to overlap mems (see "Exclusive Cpusets" on page 62)
- `notify_on_release` - Invokes `/sbin/cpuset_release_agent` when cpuset released (see "Notify on Release Flag" on page 62)
- `memory_migrate` - Causes memory pages to migrate to new mems (see "Memory Migration" on page 66)
- `memory_spread_page` - Causes kernel buffer (page) cache to spread over cpuset (see "Memory Spread" on page 65)
- `memory_spread_slab` - Causes kernel file I/O data (directory and inode slab caches) to spread over cpuset (see "Memory Spread" on page 65)

`cpuset_get_sopt`

```
const char *cpuset_get_sopt(const struct cpuset *cp, const char
*optionname);
```

Queries the value of string option `optionname` in cpuset `cp`.

Returns pointer to value of `optionname` is recognized, else returns NULL. String values in an uninitialized cpuset have value NULL.

This is an [optional] function. Use `cpuset_function()` to invoke it, as follows:

```
/* fp has function signature of pointer to cpuset_get_sopt() */
int (*fp)(struct cpuset *cp, const char *optionname);
fp = cpuset_function("cpuset_get_sopt");
if (fp) {
    fp( ... );
} else {
    puts ("cpuset_get_sopt not supported");
}
```

`cpuset_localcpus`

```
int cpuset_localcpus(const struct bitmask *mems, struct bitmask
*cpus);
```

Queries the CPUs local to specified memory nodes `mems`, by writing them to the bitmask `cpus`.

Returns 0 on success, -1 on error, setting `errno`.

cpuset_localmems

```
int cpuset_localmems(const struct bitmask *cpus, struct bitmask *mems);
```

Queries the memory nodes local to specified CPUs `cpus`, by writing them to the bitmask `mems`.

Returns 0 on success, -1 on error, setting `errno`.

cpuset_cpumemdist

```
unsigned int cpuset_cpumemdist(int cpu, int mem);
```

Distance between hardware CPU `cpu` and memory node `mem`, on a scale which has the closest distance of a CPU to its local memory valued at ten (10), and other distances more or less proportional. Distance is a rough metric of the bandwidth and delay combined, where a higher distance means lower bandwidth and longer delays.

If either `cpu` or `mem` is not known to the current system, or if any internal error occurs while evaluating this distance, or if a node has no CPUs nor memory (I/O only), then the distance returned is `UCHAR_MAX` (from `limits.h`).

These distances are obtained from the systems ACPI SLIT table, and should conform to: System Locality Information Table Interface Specification Version 1.0, July 26, 2003

This is an [optional] function. Use `cpuset_function()` to invoke it.

cpuset_cpu2node

```
int cpuset_cpu2node(int cpu);
```

Returns number of memory node closest to CPU `cpu`. For NUMA architectures (as of this writing), this commonly would be the number of the node on which `cpu` is located. If an architecture did not have memory on the same node as a CPU, it would be the node number of the memory node closest to or preferred by that `cpu`.

cpuset_create

```
int cpuset_create(const char *cpusetpath, const struct *cp);
```

Creates a cpubset at the specified `cpusetpath`, as described in the provided `struct cpubset *cp` structure. The parent cpubset of that pathname must already exist. The parameter `cp` refers to a handle obtained from a `cpuset_alloc()` call. If the parameter `cpusetpath` starts with a slash (/) character, this a path relative to `/dev/cpubset`, otherwise, it is relative to the current tasks cpubset.

The default behaviour of the `libcpubset` routine `cpuset_create()` has changed in the following way. In previous versions of SGI ProPack, the cpubset attributes `memory_spread_page`, `memory_spread_slab`, and `notify_on_release` in the newly created cpubset would default to the value zero (0) for off, regardless of the setting of these attributes in the parent cpubset. In this version of SGI ProPack, these attributes are inherited from the parent cpubset, unless explicitly set otherwise in the cpubset creation code.

Returns 0 on success, else -1 on error, setting `errno`.

This routine can fail if `malloc(3)` fails. See the `malloc(3)` man page for possible values of `errno` (`ENOMEM` being the most likely).

cpuset_delete

```
int cpubset_delete(const char *cpusetpath);
```

Deletes a cpubset at the specified `cpusetpath`. The cpubset of that pathname must already exist, be empty (no child cpubsets) and be unused (no using tasks).

If the parameter `cpusetpath` starts with a slash (/) character, this a path relative to `/dev/cpubset`, otherwise, it is relative to the current tasks cpubset.

Returns 0 on success, else -1 on error, setting `errno`.

cpuset_query

```
int cpubset_query(struct cpubset *cp, const char *cpusetpath);
```

Set `struct cpubset` structure to settings of cpubset at specified path `cpusetpath`. `struct cpubset *cp` must have been returned by a previous `cpuset_alloc()` call. Any previous settings of `cp` are lost.

If the parameter `cpusetpath` starts with a slash (/) character, this a path relative to `/dev/cpuset`, otherwise, it is relative to the current tasks `cpuset`.

Returns 0 on success, or -1 on error, setting `errno`. Errors include `cpusetpath` not referencing a valid `cpuset` path relative to `/dev/cpuset`, or the current task lacking permission to query that `cpuset`.

`cpuset_modify`

```
int cpuset_modify(const char *cpusetpath, const struct *cp);
```

Modify the `cpuset` at the specified `cpusetpath`, as described in the provided `struct cpuset *cp`. The `cpuset` at that pathname must already exist. The parameter `cp` refers to a handle obtained from a `cpuset_alloc()` call.

If the parameter `cpusetpath` starts with a slash (/) character, this a path relative to `/dev/cpuset`, otherwise, it is relative to the current tasks `cpuset`.

Returns 0 on success, else -1 on error, setting `errno`.

`cpuset_getcpusetpath`

```
char *cpuset_getcpusetpath(pid_t pid, char *buf, size_t size);
```

The `cpuset_getcpusetpath()` function copies an absolute pathname of the `cpuset` to which task of process ID `pid` is attached, to the array pointed to by `buf`, which is of length `size`. Use `pid == 0` for the current process.

The provided path is relative to the `cpuset` file system mount point.

If the `cpuset` path name would require a buffer longer than `size` elements, `NULL` is returned, and `errno` is set to `ERANGE` an application should check for this error, and allocate a larger buffer if necessary.

Returns `NULL` on failure with `errno` set accordingly, and `buf` on success. The contents of `buf` are undefined on error.

ERRORS are, as follows:

<code>EACCES</code>	Permission to read or search a component of the file name was denied.
<code>EFAULT</code>	<code>buf</code> points to a bad address.

ESRCH	The pid does not exist.
E2BIG	Larger buffer needed.
ENOSYS	Kernel does not support cpusets.

cpuset_cpusetofpid

```
int cpuset_cpusetofpid(struct cpuset *cp, int pid);
```

Set struct cpuset to settings of cpuset to which specified task pid is attached. struct cpuset *cp must have been returned by a previous cpuset_alloc() call. Any previous settings of cp are lost.

Returns 0 on success, or -1 on error, setting errno.

ERRORS are, as follows:

EACCES	Permission to read or search a component of the file name was denied.
EFAULT	buf points to a bad address.
ESRCH	The pid does not exist.
ERANGE	Larger buffer needed.
ENOSYS	Kernel does not support cpusets.

cpuset_cpusetofpid

```
int cpuset_cpusetofpid(struct cpuset *cp, int pid);
```

Sets the struct cpuset structure to settings of cpuset to which specified task pid is attached. struct cpuset *cp must have been returned by a previous cpuset_alloc() call. Any previous settings of cp are lost.

Returns 0 on success, or -1 on error, setting errno.

ERRORS are, as follows:

EACCES	Permission to read or search a component of the file name was denied.
EFAULT	buf points to a bad address.
ESRCH	The pid does not exist.

ERANGE	Larger buffer needed.
ENOSYS	Kernel does not support cpusets.

cpuset_mountpoint

```
const char *cpuset_mountpoint();
```

Returns the filesystem path at which the cpuset file system is mounted. The current implementation of this routine returns `/dev/cpuset`, or the string `[cpuset filesystem not mounted]` if the cpuset file system is not mounted, or the string `[cpuset filesystem not supported]` if the system does not support cpusets.

In general, if the first character of the return string is a slash (`/`), the result is the mount point of the cpuset file system; otherwise, the result is an error message string.

This is an [optional] function. Use `cpuset_function` to invoke it.

cpuset_collides_exclusive

```
int cpuset_collides_exclusive(const char *cpusetpath, const struct *cp);
```

Returns true (1) if cpuset `cp` would collide with any sibling of the cpuset at `cpusetpath` due to overlap of `cpu_exclusive` cpus or `mem_exclusive` mems. Return false (0) if no collision, or for any error.

The `cpuset_create` function fails with `errno == EINVAL` if the requested cpuset would overlap with any sibling, where either one is `cpu_exclusive` or `mem_exclusive`. This is a common, and not obvious error. `cpuset_collides_exclusive()` checks for this particular case, so that code creating cpusets can better identify the situation, perhaps to issue a more informative error message.

Can also be used to diagnose `cpuset_modify` failures. This routine ignores any existing cpuset with the same path as the given `cpusetpath`, and only looks for exclusive collisions with sibling cpusets of that path.

In case of any error, returns (0) – does not collide. Presumably, any actual attempt to create or modify a cpuset will encounter the same error, and report it usefully.

This routine is not particularly efficient; most likely code creating or modifying a cpuset will want to try the operation first, and then if that fails with `errno EINVAL`, perhaps call this routine to determine if an exclusive cpuset collision caused the error.

This is an [optional] function. Use `cpuset_function` to invoke it.

cpuset_init_pidlist

```
struct cpuset_pidlist *cpuset_init_pidlist(const char
*cpusetpath, int recursiveflag);
```

Initializes and returns a list of tasks (pids) attached to cpuset `cpusetpath`. If `recursiveflag` is zero, include only the tasks directly in that cpuset, otherwise, include all tasks in that cpuset or any descendant thereof.

Beware that tasks can come and go from a cpuset, after this call is made.

If the parameter `cpusetpath` starts with a slash (/) character, then this a path relative to `/dev/cpuset`, otherwise, it is relative to the current tasks cpuset.

On error, return NULL and set `errno`.

cpuset_pidlist_length

```
int cpuset_pidlist_length(const struct cpuset_pidlist *pl);
```

Returns the number of elements in `cpuset_pidlist pl`.

cpuset_get_pidlist

```
pid_t cpuset_get_pidlist(const struct cpuset_pidlist *pl, int
i);
```

Return the *i*'th element of a `cpuset_pidlist`. The elements of a `cpuset_pidlist` of length *N* are numbered 0 through *N*-1. Return `(pid_t)-1` for any other index *i*.

cpuset_free_pidlist

```
void cpuset_freepidlist(struct cpuset_pidlist *pl);
```

Deallocates a list of attached pids

cpuset_move

```
int cpuset_move(pid_t p, const char *cpusetpath);
```

Moves the task whose process ID is `p` to cpuset `cpusetpath`.

If `pid` is zero, then the current task is moved. If the parameter `cpusetpath` starts with a slash (/) character, this a path relative to `/dev/cpuset`, otherwise, it is relative to the current tasks cpuset.

Returns 0 on success, else -1 on error, setting `errno`.

cpuset_move_all

```
int cpuset_move_all(struct cpuset_pid_list *pl, const char *cpusetpath);
```

Moves all tasks in list `pl` to cpuset `cpusetpath`.

If the parameter `cpusetpath` starts with a slash (/) character, this a path relative to `/dev/cpuset`, otherwise, it is relative to the current tasks cpuset.

Returns 0 on success, else -1 on error, setting `errno`.

The `cpuset_move_all()` routine now returns an error if it was unable to move all the tasks requested.

cpuset_move_cpuset_tasks

```
int cpuset_move_cpuset_tasks(const char *fromrelpath, const char *torelpath);
```

Move all tasks in cpuset `fromrelpath` to cpuset `torelpath`. This may race with tasks being added to or forking into `fromrelpath`. Loop repeatedly, reading the task's file of cpuset `fromrelpath` and writing any task PIDs found there to the task's file of cpuset `torelpath`, up to ten attempts, or until the task's file of cpuset `fromrelpath` is empty, or until the cpuset `fromrelpath` is no longer present.

Returns 0 with `errno == 0` if able to empty the task's file of cpuset `fromrelpath`. Of course, it is still possible that some independent task could add another task to

`cpuset fromrelpath` at the same time that such a successful result is being returned. Therefore, there can be no guarantee that a successful return means that `fromrelpath` is still empty of tasks.

The `cpuset fromrelpath` might disappear during this operation, perhaps because it has `notify_on_release` flag set and was automatically removed as soon as its last task was detached from it. Consider a missing `fromrelpath` to be a successful move.

If called with `fromrelpath` and `torelpath` pathnames that evaluate to the same `cpuset`, then treat this as if `cpuset_reattach()` was called, rebinding each task in this `cpuset` one time, and return success or failure depending on the return of that `cpuset_reattach()` call.

On failure, returns -1, setting `errno`.

ERRORS are, as follows:

- `EACCES` search permission denied on intervening directory
- `ENOTEMPTY` tasks remain after multiple attempts to move them
- `EMFILE` too many open files
- `ENODEV` `/dev/cpuset` not mounted
- `ENOENT` component of `cpuset` path does not exist
- `ENOMEM` out of memory
- `ENOSYS` kernel does not support cpusets
- `ENOTDIR` component of `cpuset` path is not a directory
- `EPERM` lacked permission to read cpusets or files therein

This is an [optional] function. Use `cpuset_function` to invoke it.

`cpuset_migrate`

```
int cpuset_migrate(pid_t pid, const char *cpusetpath);
```

Migrates the task whose process ID is `p` to `cpuset cpusetpath`, moving its currently allocated memory to nodes in that `cpuset`, if not already there. If `pid` is zero, then the current task is migrated.

If the parameter `cpusetpath` starts with a slash (/) character, this a path relative to `/dev/cpuset`, otherwise, it is relative to the current tasks `cpuset`. For more information, see "Memory Migration" on page 66.

Returns 0 on success, else -1 on error, setting `errno`.

This is an [optional] function. Use `cpuset_function()` to invoke it.

`cpuset_migrate_all`

```
int cpuset_migrate_all(struct cpuset_pid_list *pl, const char
*cpusetpath);
```

Moves all tasks in list `pl` to `cpuset cpusetpath`, moving their currently allocated memory to nodes in that `cpuset`, if not already there.

If the parameter `cpusetpath` starts with a slash (/) character, this a path relative to `/dev/cpuset`, otherwise, it is relative to the current tasks `cpuset`.

This is an [optional] function. Use `cpuset_function()` to invoke it.

`cpuset_reattach`

```
int cpuset_reattach(const char *cpusetpath);
```

Reattaches all tasks in `cpuset cpusetpath` to itself. This additional step is necessary anytime that the `cpus` of a `cpuset` have been changed, in order to rebind the `cpus_allowed` of each task in the `cpuset` to the new value. This routine writes the `pid` of each task currently attached to the named `cpuset` to the tasks file of that `cpuset`. If additional tasks are being spawned too rapidly into the `cpuset` at the same time, there is an unavoidable race condition, and some tasks may be missed.

If the parameter `cpusetpath` starts with a slash (/) character, this a path relative to `/dev/cpuset`, otherwise, it is relative to the current tasks `cpuset`. Returns 0 on success, else -1 on error, setting `errno`.

`cpuset_open_memory_pressure`

```
int cpuset_open_memory_pressure(const char *cpusetpath);
```

Opens a file descriptor from which to read the `memory_pressure` of the cuset `cpusetpath`.

If the `cpusetpath` parameter starts with a slash (/) character, this a path relative to `/dev/cpuset`; otherwise, it is relative to the current tasks cuset.

By default, computation by the kernel of `memory_pressure` is disabled. Set the `memory_pressure_enabled` flag in the top cuset to enable it.

For more information, see "Memory Pressure of a Cuset" on page 63.

This is an [optional] function. Use the `cpuset_function` to invoke it.

`cpuset_read_memory_pressure`

```
int cpuset_read_memory_pressure(int fd);
```

Reads and return the current `memory_pressure` of the cuset for which file descriptor `fd` was opened using the `cpuset_open_memory_pressure` function.

Uses the system call `pread(2)`. On success, returns a non-negative number, as described in "Memory Pressure of a Cuset" on page 63. On failure, returns -1 and sets `errno`.

By default, computation by the kernel of `memory_pressure` is disabled. Set the `memory_pressure_enabled` flag in the top cuset to enable it.

For more information, see "Memory Pressure of a Cuset" on page 63.

This is an [optional] function. Use the `cpuset_function` to invoke it.

`cpuset_close_memory_pressure`

```
void cpuset_close_memory_pressure(int fd);
```

Closes the file descriptor `fd` which was opened using the `cpuset_open_memory_pressure` function.

If `fd` is not a valid open file descriptor, this call does nothing. No error is returned in any case.

By default, computation by the kernel of `memory_pressure` is disabled. Set the `memory_pressure_enabled` flag in the top cuset to enable it.

For more information, see "Memory Pressure of a Cpuset" on page 63.

This is an [optional] function. Use the `cpuset_function` to invoke it.

`cpuset_c_rel_to_sys_cpu`

```
int cpuset_c_rel_to_sys_cpu(const struct cpuset *cp, int cpu);
```

Returns the system-wide CPU number that is used by the `cpu`-th CPU of the specified `cpuset cp`. Returns result of `cpuset_cpus_nbits()` if `cpu` is not in the range `[0, bitmask_weight(cpuset_getcpus(cp))]`.

`cpuset_c_sys_to_rel_cpu`

```
int cpuset_c_sys_to_rel_cpu(const struct cpuset *cp, int cpu);
```

Returns the `cpu`-th CPU of the specified `cpuset cp` that is used by the system-wide CPU number. Returns result of `cpuset_cpus_nbits()` if `bitmask_isbitset(cpuset_getcpus(cp), cpu)` is false.

`cpuset_c_rel_to_sys_mem`

```
int cpuset_c_rel_to_sys_mem(const struct cpuset *cp, int mem);
```

Returns the system-wide memory node number that is used by the `mem`-th memory node of the specified `cpuset cp`. Returns result of `cpuset_mems_nbits()` if `mem` is not in the range `[0, bitmask_weight(cpuset_getmems(cp))]`. Note that this is a left closed, right open interval. The set of points in the interval `[a, b)` is the set of all points `x` such that `a <= x < b`.

`cpuset_c_sys_to_rel_mem`

```
int cpuset_c_sys_to_rel_mem(const struct cpuset *cp, int mem);
```

Returns the `mem`-th memory node of the specified `cpuset cp` that is used by the system-wide memory node number. Returns result of `cpuset_mems_nbits()` if `bitmask_isbitset(cpuset_getmems(cp), mem)` is false.

cpuset_p_rel_to_sys_cpu

```
int cpuset_p_rel_to_sys_cpu(pid_t pid, int cpu);
```

Returns the system-wide CPU number that is used by the `cpu`-th CPU of the cpuset to which task `pid` is attached. Returns result of `cpuset_cpus_nbits()` if that cpuset does not encompass that relative cpu number.

cpuset_p_sys_to_rel_cpu

```
int cpuset_p_sys_to_rel_cpu(pid_t pid, int cpu);
```

Returns the `cpu`-th CPU of the cpuset to which task `pid` is attached that is used by the system-wide CPU number. Returns result of `cpuset_cpus_nbits()` if that cpuset does not encompass that system-wide cpu number.

cpuset_p_rel_to_sys_mem

```
int cpuset_p_rel_to_sys_mem(pid_t pid, int mem);
```

Returns the system-wide memory node number that is used by the `mem`-th memory node of the cpuset to which task `pid` is attached. Returns result of `cpuset_mems_nbits()` if that cpuset does not encompass that relative memory node number.

cpuset_p_sys_to_rel_mem

```
int cpuset_p_sys_to_rel_mem(pid_t pid, int mem);
```

Returns the `mem`-th memory node of the cpuset to which task `pid` is attached that is used by the system-wide memory node number. Returns result of `cpuset_mems_nbits()` if that cpuset does not encompass that system-wide memory node.

cpuset_get_placement

```
cpuset_get_placement(pid) - [optional Return current placement  
of task pid]
```

This function returns an opaque struct placement * pointer. The results of calling `cpuset_get_placement` twice at different points in a program can be compared using `cpuset_equal_placement` to determine if the specified task has had its cpuset CPU and memory placement modified between those two `cpuset_get_placement` calls.

When finished with a struct placement * pointer, free it by calling `cpuset_free_placement`.

This is an [optional] function. Use the `cpuset_function` to invoke it.

cpuset_equal_placement

```
cpuset_equal_placement(plc1, plc2) - [optional] True if two  
placements equal
```

This function compares two struct placement * pointers, returned by two separate calls to `cpuset_get_placement`. This is done to determine if the specified task has had its cpuset CPU and memory placement modified between those two `cpuset_get_placement` calls.

When finished with a struct placement * pointer, free it by calling the `cpuset_free_placement` function.

Two struct placement * pointers will compare equal if they have the same CPU placement `cpus`, the same memory placement `mems`, and the same cpuset path.

This is an [optional] function. Use the `cpuset_function` to invoke it.

cpuset_free_placement

```
cpuset_free_placement(plc) - [optional] Free placement
```

Use this routine to free a struct placement * pointer returned by a previous call to the `cpuset_get_placement` function.

This is an [optional] function. Use the `cpuset_function` to invoke it.

cpuset_cpubind

```
int cpuset_cpubind(int cpu);
```

Binds the scheduling of the current task to CPU `cpu`, using the `sched_setaffinity(2)` system call.

Fails with a return of `-1`, and `errno` set to `EINVAL`, if `cpu` is not allowed in the current tasks `cpuset`.

The following code will bind the scheduling of a thread to the `n`-th CPU of the current `cpuset`:

```
/*
 * Confine current task to only run on the n-th CPU
 * of its current cpuset. If in a cpuset of N CPUs,
 * valid values for n are 0 .. N-1.
 */
cpuset_cpupbind(cpuset_p_rel_to_sys_cpu(0, n));
```

cpuset_latestcpu

```
int cpuset_latestcpu(pid_t pid);
```

Returns the most recent CPU on which the specified task `pid` executed. If `pid` is `0`, examine current task.

The `cpuset_latestcpu()` call returns the number of the CPU on which the specified task `pid` most recently executed. If a process can be scheduled on two or more CPUs, the results of `cpuset_lastcpu()` may become invalid even before the query returns to the invoking user code.

The last used CPU is visible for a given `pid` as field #39 (starting with #1) in the file `/proc/pid/stat`. Currently, this file has 41 fields, so it is the 3rd to the last field.

cpuset_membind

```
int cpuset_membind(int mem);
```

Binds the memory allocation of the current task to memory node `mem`, using the `set_mempolicy(2)` system call with a policy of `MPOL_BIND`.

Fails with a return of `-1`, and `errno` set to `EINVAL`, if `mem` is not allowed in the current tasks `cpuset`.

The following code will bind the memory allocation of a thread to the n-th memory node of the current cpuset:

```
/*
 * Confine current task to only allocate memory on
 * n-th Node of its current cpuset.  If in a cpuset
 * of N Memory Nodes, valid values for n are 0 .. N-1.
 */
cpuset_membind(cpuset_p_rel_to_sys_mem(0, n));
```

cpuset_nuke

```
int cpuset_nuke(const char *cpusetpath, unsigned int seconds);
```

Remove a cpuset, including killing tasks in it, and removing any descendent cpusets and killing their tasks.

Tasks can take a long time (minutes on some configurations) to exit. Loop up to seconds seconds, trying to kill them.

The following steps are taken to remove a cpuset:

- First, kills all the PIDs, looping until there are no more PIDs in this cpuset or below or until the `seconds` timeout limit is exceeded.
- Second, remove that cpuset, and any child cpusets it has, starting from the lowest level leaf node cpusets and working back upwards.
- Third, if by this point the original cpuset is gone, return success.

If the timeout is exceeded, and tasks still exist, fail with `errno == ETIME`.

This routine sleeps a variable amount of time. After the first attempt to kill all the tasks in the cpuset or its descendents, it sleeps one second, the next time it sleeps two seconds, increasing one second each loop up to a maximum of ten seconds. If more loops past ten seconds are required to kill all the tasks, it sleeps ten seconds each subsequent loop. In any case, before the last loop, it sleeps however many seconds remain of the original timeout seconds requested. The total time of all sleeps will be no more than the requested seconds.

If the cpuset started out empty of any tasks, or if the passed in seconds was zero, this routine will return quickly, having not slept at all. Otherwise, this routine will at a minimum send a `SIGKILL` signal to all the tasks in this cpuset subtree, then sleep

one second, before looking to see if any tasks remain. If tasks remain in the cpuset subtree, and a longer seconds timeout was requested (more than one), it will continue to kill remaining tasks and sleep, in a loop, for as long as time and tasks remain.

`cpuset_export`

```
int cpuset_export(const struct cpuset *cp, char *buf, int
buflen);
```

Writes the settings of cpuset `cp` to file. If no file exists at the path specified by `file`, create one. If a file already exists there, overwrite it.

Returns -1 and sets `errno` on error. Upon successful return, returns the number of characters printed (not including the trailing '0' used to end output to strings). The function `cpuset_export` does not write more than `size` bytes (including the trailing '0'). If the output was truncated due to this limit, the return value is the number of characters (not including the trailing '0') which would have been written to the final string if enough space had been available. Thus, a return value of `size` or more means that the output was truncated.

For details of the format required for exported cpuset file, see "Cpuset Text Format" on page 82.

`cpuset_import`

```
int cpuset_import(struct cpuset *cp, const char *file, int
*errlinenum_ptr, char *errmsg_bufptr, int errmsg_buflen);
```

Reads the settings of cpuset `cp` from file. If no file exists at the path specified by `file`, create one. If a file already exists there, overwrite it.

`struct cpuset *cp` must have been returned by a previous `cpuset_alloc()` call. Any previous settings of `cp` are lost.

Returns 0 on success, or -1 on error, setting `errno`. Errors include file not referencing a readable file.

If parsing errors are encountered reading the file and if `errlinenum_ptr` is not NULL, the number of the first line (numbers start with one) with an error is written to `*errlinenum_ptr`. If an error occurs on the open and `errlinenum_ptr` is not NULL, zero is written to `*errlinenum_ptr`.

If parsing errors are encountered reading the file and if `errmsg_bufptr` is not NULL, it is presumed to point to a character buffer of at least `errmsg_buflen` characters and a nul-terminated error message is written to `*errmsg_bufptr`, providing a human readable error message explaining the error message in more detail. Currently, the possible error messages are, as follows:

- "Token 'CPU' requires list"
- "Token 'MEM' requires list"
- "Invalid list format: %s"
- "Unrecognized token: %s"
- "Insufficient memory"

For details of the format required for imported cpuset file, see "Cpuset Text Format" on page 82.

`cpuset_function`

```
cpuset_function(const char *function_name);
```

Returns pointer to the named `libcpuset.so` function, or NULL. For base functions that are in all implementations of `libcpuset`, there is no particular value in using `cpuset_function()` to obtain a pointer to the function dynamically. But for [optional] cpuset functions, the use of `cpuset_function()` enables dynamically adapting to runtime environments that may or may not support that function.

`cpuset_version`

```
int cpuset_version();
```

Version (simple integer) of the cpuset library (`libcpuset`). The version number returned by `cpuset_version()` is incremented anytime that any changes or additions are made to its API or behaviour. Other mechanisms are provided to maintain full upward compatibility with this libraries API. This `cpuset_version()` call is intended to provide a fallback mechanism in case an application needs to distinguish between two previous versions of this library.

This is an [optional] function. Use `cpuset_function` to invoke it.

Functions to Traverse a Cpuset Hierarchy

The functions described in this section are used to transverse a cpuset hierarchy.

`cpuset_fts_open`

```
struct cpuset_fts_tree *cpuset_fts_open(const char *cpusetpath);
```

Opens a cpuset hierarchy. Returns a pointer to a `cpuset_fts_tree` structure, which can be used to traverse all cpusets below the specified cpuset `cpusetpath`.

If the parameter `cpusetpath` starts with a slash (/) character, then this path is relative to `/dev/cpuset`, otherwise, it is relative to the current tasks cpuset.

The `cpuset_fts_open` routine is implemented internally using the `fts(3)` library routines for traversing a file hierarchy. The entire cpuset subtree below `cpusetpath` is traversed as part of the `cpuset_fts_open()` call, and all cpuset state and directory `stat` information is captured at that time. The other `cpuset_fts_*` routines just access this captured state. Any changes to the traversed cpusets made after the return of the `cpuset_fts_open()` call will not be visible via the returned `cpuset_fts_tree` structure.

Internally, the `fts(3)` options `FTS_NOCHDIR` and `FTS_XDEV` are used, to avoid changing the invoking tasks current directory, and to avoid descending into any other file systems mounted below `/dev/cpuset`. The order in which cpusets will be returned by the `cpuset_fts_read` routine corresponds to the `fts` pre-order (`FTS_D`) visitation order. The internal `fts` scan by `cpuset_fts_open` ignores the post-order (`FTS_DP`) results.

Because the `cpuset_fts_open()` call collects all the information at once from an entire cpuset subtree, a simple error return would not provide sufficient information as to what failed, and on what cpuset in the subtree. So, except for `malloc(3)` failures, errors are captured in the list of entries.

See `cpuset_fts_get_info` for details of the information field.

This is an [optional] function. Use `cpuset_function` to invoke it.

`cpuset_fts_read`

```
const struct cpuset_fts_entry *cpuset_fts_read(struct cpuset_fts_tree *cs_tree);
```

Returns next `cs_entry` in `cpuset_fts_tree cs_tree` obtained from an `cpuset_fts_open()` call. One `cs_entry` is returned for each cpuset directory that

was found in the subtree scanned by the `cpuset_fts_open()` call. Use the `info` field obtained from a `cpuset_fts_get_info()` call to determine which fields of a particular `cs_entry` are valid, and which fields contain error information or are not valid.

This is an [optional] function. Use `cpuset_function` to invoke it.

`cpuset_fts_reverse`

```
void cpuset_fts_reverse(struct cpuset_fts_tree *cs_tree);
```

Reverse order of `cs_entry`'s in the `cpuset_fts_tree` `cs_tree` obtained from a `cpuset_fts_open()` call.

An open `cpuset_fts_tree` stores a list of `cs_entry` cpuset entries, in pre-order, meaning that a series of `cpuset_fts_read()` calls will always return a parent cpuset before any of its child cpusets. Following a `cpuset_fts_reverse()` call, the order of cpuset entries is reversed, putting it in post-order, so that a series of `cpuset_fts_read()` calls will always return any children cpusets before their parent cpuset. A second `cpuset_fts_reverse()` call would put the list back in pre-order again.

To avoid exposing confusing inner details of the implementation across the API, a `cpuset_fts_rewind()` call is always automatically performed on a `cpuset_fts_tree` whenever `cpuset_fts_reverse()` is called on it.

This is an [optional] function. Use `cpuset_function` to invoke it.

`cpuset_fts_rewind`

```
void cpuset_fts_rewind(struct cpuset_fts_tree *cs_tree);
```

Rewind a cpuset tree `cs_tree` obtained from a `cpuset_fts_open()` call, so that subsequent `cpuset_fts_read()` calls start from the beginning again.

This is an [optional] function. Use `cpuset_function` to invoke it.

`cpuset_fts_get_path`

```
const char *cpuset_fts_get_path(const struct cpuset_fts_entry *cs_entry);
```

Return the cpuset path, relative to `/dev/cpuset`, as null-terminated string, of a `cs_entry` obtained from a `cpuset_fts_read()` call.

The results of this call are valid for all `cs_entry`'s returned from `cpuset_fts_read()` calls, regardless of the value returned by `cpuset_fts_get_info()` for that `cs_entry`.

This is an [optional] function. Use `cpuset_function` to invoke it.

`cpuset_fts_get_stat`

```
const struct stat *cpuset_fts_get_stat(const struct cpuset_fts_entry *cs_entry);
```

Return pointer to `stat(2)` information about the `cpuset` directory of a `cs_entry` obtained from a `cpuset_fts_read()` call.

The results of this call are valid for all `cs_entry`'s returned from `cpuset_fts_read()` calls, regardless of the value returned by `cpuset_fts_get_info()` for that `cs_entry`, except in the cases that:

- The information field returned by `cpuset_fts_get_info` contains `CPUSET_FTS_ERR_DNR`, in which case, a directory in the path to the `cpuset` could not be read and this call will return a NULL pointer.
- The information field returned by `cpuset_fts_get_info` contains `CPUSET_FTS_ERR_STAT`, in which case a `stat(2)` failed on this `cpuset` directory and this call will return a pointer to a `struct stat` containing all zeros.

This is an [optional] function. Use `cpuset_function` to invoke it.

`cpuset_fts_get_cpuset`

```
const struct cpuset *cpuset_fts_get_cpuset(const struct cpuset_fts_entry *cs_entry);
```

Return the `struct cpuset` pointer of a `cs_entry` obtained from a `cpuset_fts_read()` call. The `struct cpuset` so referenced describes the `cpuset` represented by one directory in the `cpuset` hierarchy, and can be used with various other calls in this library.

The results of this call are only valid for a `cs_entry` if the `cpuset_fts_get_info()` call returns `CPUSET_FTS_CPUSSET` for the information field of a `cs_entry`. If the info field contained `CPUSET_FTS_ERR_CPUSSET`, then `cpuset_fts_get_cpuset` returns a pointer to a `struct cpuset` that is all zeros. If the information field contains any other `CPUSET_FTS_ERR_*` value, then `cpuset_fts_get_cpuset` returns a NULL pointer.

This is an [optional] function. Use `cpuset_function` to invoke it.

cpuset_fts_get_errno

```
int cpuset_fts_get_errno(const struct cpuset_fts_entry *cs_entry);
```

Returns the error field of a `cs_entry` obtained from a `cpuset_fts_read()` call.

If this information field has one of the following `CPUSET_FTS_ERR_*` values, it indicates which operation failed, the error field (returned by `cpuset_fts_get_errno`) captures the failing `errno` value for that operation, the `path` field (returned by `cpuset_fts_get_path`) indicates which cpuset failed, and some of the other entry fields may not be valid, depending on the value. If an entry has the value `CPUSET_FTS_CPUSET` for its information field, then the error field will have the value 0, and the other fields will be contain valid information about that cpuset.

Information field values are, as follows:

```
CPUSET_FTS_ERR_DNR = 0:  
Error - couldn't read directory  
CPUSET_FTS_ERR_STAT = 1:  
Error - couldn't stat directory  
CPUSET_FTS_ERR_CPUSET = 2:  
Error - cpuset_query failed  
CPUSET_FTS_CPUSET = 3:  
Valid cpuset
```

The above information field values are defined using an anonymous enum in the `cpuset.h` header file. If it necessary to maintain source code compatibility with earlier versions of the `cpuset.h` header file lacking the above `CPUSET_FTS_*` values, one can conditionally check that the C preprocessor symbol `CPUSET_FTS_INFO_VALUES_DEFINED` is not defined and provide alternative coding for that case.

This is an [optional] function. Use `cpuset_function` to invoke it.

cpuset_fts_close

```
void cpuset_fts_close(struct cpuset_fts_tree *cs_tree);
```

Close a `cs_tree` obtained from a `cpuset_fts_open()` call, freeing any internally allocated memory for that `cs_tree`.

This is an [optional] function. Use `cpuset_function` to invoke it.

Index

A

Array Services, 3

- accessing an array, 8
- array configuration database, 1, 3
- array daemon, 3
- array name, 8
- array session handle, 1, 20
- arrayconfig_tempo command, 6

ASH

- See "array session handle", 1

authentication key, 14

commands, 3

- ainfo, 3, 8, 13, 14
- array, 3, 14
- arshell, 3, 14

common command options, 14

common environment variables, 16

concepts

- array session, 13
- array session handle, 13

ASH

- See "array session handle", 13

configuring, 6

finding basic usage information, 8

global process namespace, 1

hostname command, 14

ibarray, 3

invoking a program, 9

- information sources, 10
- ordinary (sequential) applications, 9
- parallel message-passing applications
 - distributed over multiple nodes , 9
- parallel message-passing applications
 - within a node, 9
- parallel shared-memory applications within
 - a node, 9

- local process management commands, 12
 - at, 12
 - batch, 12
 - crontab, 12
 - intro, 12
 - kill, 12
 - nice, 12
 - ps, 12
 - top, 12
- logging into an array, 8
- managing local processes, 11
- monitoring processes and system usage, 11
- names of arrays and nodes, 13
- overview, 1
- release notes, 2
- scheduling and killing local processes, 11
- security considerations, 26
- specifying a single node, 16
- using an array, 7
- using array services commands, 12

arrayconfig_tempo command, 6

C

Configuring Array Services, 6

Cpuset Facility

- advantages, 47
- boot cpuset, 76
 - creating, 77
 - /etc/bootcpuset.conf file, 78
- command line utility, 72

cpuset

- definition, 47
- determine if cpusets are installed, 51
- errno settings, 88
- modifying CPUs and kernel processing, 83

- system error messages, 88
 - cpuset permissions, 67
 - cpuset text format, 82
 - directories, 56
 - overview, 47
 - programming model, 86
 - scheduling and memory allocation, 68
 - systems calls
 - mbind, 48
 - sched_setaffinity, 48
 - set_mempolicy, 48
 - using cpusets at shell prompt, 70
 - create a cpuset, 70
 - remove a cpuset, 71
 - Cpuset Facility on SGI Performance Suite
 - Cpuset library functions
 - cpuset_alloc, 102
 - cpuset_c_rel_to_sys_cpu, 118
 - cpuset_c_rel_to_sys_mem, 118
 - cpuset_c_sys_to_rel_cpu, 118
 - cpuset_c_sys_to_rel_mem, 118
 - cpuset_close_memory_pressure, 117
 - cpuset_collides_exclusive, 112
 - cpuset_cpu2node, 108
 - cpuset_cpubind, 121
 - cpuset_cpus_nbits, 103
 - cpuset_cpus_weight, 106
 - cpuset_cpusetofpid, 111
 - cpuset_create, 109
 - cpuset_delete, 109
 - cpuset_equal_placement, 120
 - cpuset_export, 123
 - cpuset_free, 103
 - cpuset_free_pidlist, 114
 - cpuset_free_placement, 120
 - cpuset_function, 124
 - cpuset_get_iopt, 106
 - cpuset_get_pidlist, 113
 - cpuset_get_placement, 120
 - cpuset_get_sopt, 107
 - cpuset_getcpus, 105
 - cpuset_getcpusetpath, 110
 - cpuset_getmems, 106
 - cpuset_import, 123
 - cpuset_init_pidlist, 113
 - cpuset_latestcpu, 121
 - cpuset_localcpus, 108
 - cpuset_localmems, 108
 - cpuset_membind, 121
 - cpuset_mems_nbits, 103
 - cpuset_mems_weight, 106
 - cpuset_migrate, 115
 - cpuset_migrate_all, 116
 - cpuset_mountpoint, 112
 - cpuset_move, 114
 - cpuset_move_all, 114
 - cpuset_open_memory_pressure, 117
 - cpuset_p_rel_to_sys_cpu, 119
 - cpuset_p_rel_to_sys_mem, 119
 - cpuset_p_sys_to_rel_cpu, 119
 - cpuset_p_sys_to_rel_mem, 119
 - cpuset_pidlist_length, 113
 - cpuset_pin, 95
 - cpuset_query, 109
 - cpuset_reattach, 116
 - cpuset_set_iopt, 104
 - cpuset_set_sopt, 105
 - cpuset_setcpus, 103
 - cpuset_setmems, 104
 - cpuset_size, 96
 - cpuset_unpin, 96
 - cpuset_where, 96
 - puset_cpumemdist, 108
 - puset_modify, 110
 - puset_read_memory_pressure, 117
 - Creating a cpuset, 70
- N
- NUMA Tools
 - Command
 - dlook, 91

R

Removing a cpuset, 71

S

Secure Array Services
certificates, 40
cert directory, 41
keys, 41
commands

gencert, 41
makecert, 42
differences with standard Array Services, 39
overview, 39
parameters
ssl_cipher_list, 44
ssl_verify_depth, 44
secure shell considerations, 44
System configuration
Configuring Array Services, 6