

Resource Management and Oracle® Solaris Zones Developer's Guide

ORACLE®

Part No: E54826
October 2017

Part No: E54826

Copyright © 2004, 2017, Oracle and/or its affiliates. All rights reserved.

This software and related documentation are provided under a license agreement containing restrictions on use and disclosure and are protected by intellectual property laws. Except as expressly permitted in your license agreement or allowed by law, you may not use, copy, reproduce, translate, broadcast, modify, license, transmit, distribute, exhibit, perform, publish, or display any part, in any form, or by any means. Reverse engineering, disassembly, or decompilation of this software, unless required by law for interoperability, is prohibited.

The information contained herein is subject to change without notice and is not warranted to be error-free. If you find any errors, please report them to us in writing.

If this is software or related documentation that is delivered to the U.S. Government or anyone licensing it on behalf of the U.S. Government, then the following notice is applicable:

U.S. GOVERNMENT END USERS: Oracle programs, including any operating system, integrated software, any programs installed on the hardware, and/or documentation, delivered to U.S. Government end users are "commercial computer software" pursuant to the applicable Federal Acquisition Regulation and agency-specific supplemental regulations. As such, use, duplication, disclosure, modification, and adaptation of the programs, including any operating system, integrated software, any programs installed on the hardware, and/or documentation, shall be subject to license terms and license restrictions applicable to the programs. No other rights are granted to the U.S. Government.

This software or hardware is developed for general use in a variety of information management applications. It is not developed or intended for use in any inherently dangerous applications, including applications that may create a risk of personal injury. If you use this software or hardware in dangerous applications, then you shall be responsible to take all appropriate fail-safe, backup, redundancy, and other measures to ensure its safe use. Oracle Corporation and its affiliates disclaim any liability for any damages caused by use of this software or hardware in dangerous applications.

Oracle and Java are registered trademarks of Oracle and/or its affiliates. Other names may be trademarks of their respective owners.

Intel and Intel Xeon are trademarks or registered trademarks of Intel Corporation. All SPARC trademarks are used under license and are trademarks or registered trademarks of SPARC International, Inc. AMD, Opteron, the AMD logo, and the AMD Opteron logo are trademarks or registered trademarks of Advanced Micro Devices. UNIX is a registered trademark of The Open Group.

This software or hardware and documentation may provide access to or information about content, products, and services from third parties. Oracle Corporation and its affiliates are not responsible for and expressly disclaim all warranties of any kind with respect to third-party content, products, and services unless otherwise set forth in an applicable agreement between you and Oracle. Oracle Corporation and its affiliates will not be responsible for any loss, costs, or damages incurred due to your access to or use of third-party content, products, or services, except as set forth in an applicable agreement between you and Oracle.

Access to Oracle Support

Oracle customers that have purchased support have access to electronic support through My Oracle Support. For information, visit <http://www.oracle.com/pls/topic/lookup?ctx=acc&id=info> or visit <http://www.oracle.com/pls/topic/lookup?ctx=acc&id=trs> if you are hearing impaired.

Référence: E54826

Copyright © 2004, 2017, Oracle et/ou ses affiliés. Tous droits réservés.

Ce logiciel et la documentation qui l'accompagne sont protégés par les lois sur la propriété intellectuelle. Ils sont concédés sous licence et soumis à des restrictions d'utilisation et de divulgation. Sauf stipulation expresse de votre contrat de licence ou de la loi, vous ne pouvez pas copier, reproduire, traduire, diffuser, modifier, accorder de licence, transmettre, distribuer, exposer, exécuter, publier ou afficher le logiciel, même partiellement, sous quelque forme et par quelque procédé que ce soit. Par ailleurs, il est interdit de procéder à toute ingénierie inverse du logiciel, de le désassembler ou de le décompiler, excepté à des fins d'interopérabilité avec des logiciels tiers ou tel que prescrit par la loi.

Les informations fournies dans ce document sont susceptibles de modification sans préavis. Par ailleurs, Oracle Corporation ne garantit pas qu'elles soient exemptes d'erreurs et vous invite, le cas échéant, à lui en faire part par écrit.

Si ce logiciel, ou la documentation qui l'accompagne, est livré sous licence au Gouvernement des Etats-Unis, ou à quiconque qui aurait souscrit la licence de ce logiciel pour le compte du Gouvernement des Etats-Unis, la notice suivante s'applique :

U.S. GOVERNMENT END USERS: Oracle programs, including any operating system, integrated software, any programs installed on the hardware, and/or documentation, delivered to U.S. Government end users are "commercial computer software" pursuant to the applicable Federal Acquisition Regulation and agency-specific supplemental regulations. As such, use, duplication, disclosure, modification, and adaptation of the programs, including any operating system, integrated software, any programs installed on the hardware, and/or documentation, shall be subject to license terms and license restrictions applicable to the programs. No other rights are granted to the U.S. Government.

Ce logiciel ou matériel a été développé pour un usage général dans le cadre d'applications de gestion des informations. Ce logiciel ou matériel n'est pas conçu ni n'est destiné à être utilisé dans des applications à risque, notamment dans des applications pouvant causer un risque de dommages corporels. Si vous utilisez ce logiciel ou ce matériel dans le cadre d'applications dangereuses, il est de votre responsabilité de prendre toutes les mesures de secours, de sauvegarde, de redondance et autres mesures nécessaires à son utilisation dans des conditions optimales de sécurité. Oracle Corporation et ses affiliés déclinent toute responsabilité quant aux dommages causés par l'utilisation de ce logiciel ou matériel pour des applications dangereuses.

Oracle et Java sont des marques déposées d'Oracle Corporation et/ou de ses affiliés. Tout autre nom mentionné peut correspondre à des marques appartenant à d'autres propriétaires qu'Oracle.

Intel et Intel Xeon sont des marques ou des marques déposées d'Intel Corporation. Toutes les marques SPARC sont utilisées sous licence et sont des marques ou des marques déposées de SPARC International, Inc. AMD, Opteron, le logo AMD et le logo AMD Opteron sont des marques ou des marques déposées d'Advanced Micro Devices. UNIX est une marque déposée de The Open Group.

Ce logiciel ou matériel et la documentation qui l'accompagne peuvent fournir des informations ou des liens donnant accès à des contenus, des produits et des services émanant de tiers. Oracle Corporation et ses affiliés déclinent toute responsabilité ou garantie expresse quant aux contenus, produits ou services émanant de tiers, sauf mention contraire stipulée dans un contrat entre vous et Oracle. En aucun cas, Oracle Corporation et ses affiliés ne sauraient être tenus pour responsables des pertes subies, des coûts occasionnés ou des dommages causés par l'accès à des contenus, produits ou services tiers, ou à leur utilisation, sauf mention contraire stipulée dans un contrat entre vous et Oracle.

Accès aux services de support Oracle

Les clients Oracle qui ont souscrit un contrat de support ont accès au support électronique via My Oracle Support. Pour plus d'informations, visitez le site <http://www.oracle.com/pls/topic/lookup?ctx=acc&id=info> ou le site <http://www.oracle.com/pls/topic/lookup?ctx=acc&id=trs> si vous êtes malentendant.

Contents

- Using This Documentation** 9

- 1 Resource Management in the Oracle Solaris Operating System** 11
 - Understanding Resource Management in the Oracle Solaris Operating System 11
 - Resource Management Workload Organization 11
 - Resource Organization 12
 - Resource Controls 13
 - Extended Accounting Facility 14
 - Writing Resource Management Applications 14

- 2 Workload Hierarchy of Projects and Tasks** 15
 - Overview of Projects and Tasks 15
 - /etc/project File 16
 - Project and Task API Functions 18
 - Code Examples for Accessing project Database Entries 19
 - Programming Issues Associated With Projects and Tasks 20

- 3 Using the C Interface to Extended Accounting** 21
 - Overview of the C Interface to Extended Accounting 21
 - Extended Accounting API Functions 21
 - exact System Calls 22
 - Operations on the exact File 22
 - Operations on exact Objects 22
 - Extended Accounting Memory Management 23
 - Extended Accounting Miscellaneous Operations 24
 - C Code Examples for Accessing exact Files 24
 - Programming Issues With exact Files 28

4 Using the Perl Interface to Extended Accounting	29
Extended Accounting Overview	29
Perl Interface to libexacct	30
Perl Interface Object Model	30
Benefits of Using the Perl Interface to libexacct	30
Perl Double-Typed Scalars	31
Perl Modules	31
Sun::Solaris::Project Module	33
Sun::Solaris::Task Module	34
Sun::Solaris::Exacct Module	35
Sun::Solaris::Exacct::Catalog Module	36
Sun::Solaris::Exacct::File Module	38
Sun::Solaris::Exacct::Object Module	40
Sun::Solaris::Exacct::Object::Item Module	41
Sun::Solaris::Exacct::Object::Group Module	42
Sun::Solaris::Exacct::Object::_Array Module	43
Perl Code Examples	43
Output From the Perl dump Method	47
5 Resource Controls in Oracle Solaris	51
Overview of Resource Controls	51
Resource Controls Flags and Actions	52
rlimit, Resource Limit	52
rctl, Resource Control	52
Resource Control Values and Privilege Levels	53
Local Actions and Local Flags	53
Global Actions and Global Flags	54
Resource Control Sets Associated With a Zone, Project, Processes, and Tasks	55
Signals Used With Resource Controls	62
Resource Controls API Functions	63
Operate on Action-Value Pairs of a Resource Control	63
Operate on Local Modifiable Values	64
Retrieve Local Read-Only Values	64
Retrieve Global Read-Only Actions	64
Resource Control Code Examples	64

Master Observing Process for Resource Controls	64
List all the Value-Action Pairs for a Specific Resource Control	66
Set <code>project.cpu-shares</code> and Add a New Value	67
Set LWP Limit Using Resource Control Blocks	68
Programming Issues Associated With Resource Controls	69
<code>zonestat</code> Utility for Monitoring Zones Resource Usage	69
6 Resource Pools in Oracle Solaris	71
Overview of Resource Pools	71
Scheduling Class	72
Dynamic Resource Pool Constraints and Objectives	72
System Properties	73
Pools Properties	74
Processor Set Properties	75
Using <code>libpool</code> to Manipulate Pool Configurations	76
Manipulate <code>psets</code>	77
Resource Pools API Functions	77
Functions for Operating on Resource Pools and Associated Elements	78
Functions for Querying Resource Pools and Associated Elements	80
Resource Pool Code Examples	83
Ascertain the Number of CPUs in the Resource Pool	83
List All Resource Pools	84
Report Pool Statistics for a Given Pool	85
Set <code>pool.comment</code> Property and Add New Property	85
Programming Issues Associated With Resource Pools	86
<code>zonestat</code> Utility for Monitoring Resource Pools in Oracle Solaris Zones	87
7 Design Considerations for Resource Management Applications in Oracle Solaris Zones	89
Oracle Solaris Zones Overview	89
IP Networking in Oracle Solaris Zones	90
About Applications in Oracle Solaris Zones	90
General Considerations When Writing Applications for Non-Global Zones	91
Specific Considerations for Oracle Solaris 10 Non-Global Zones	93
Specific Considerations for Shared-IP Non-Global Zones	93
Packaging Considerations in <code>solaris</code> Zones	94

API for Zones Monitoring Statistics	94
Monitoring Zone File System Activity	95
Oracle Solaris 10 Zones	96
Oracle Solaris Kernel Zones	96
8 Project Configuration Examples	99
/etc/project Project File	99
Define Two Projects	99
Configure Resource Controls	100
Configure Resource Pools	100
Configure FSS project.cpu-shares for a Project	100
Configure Five Applications with Different Characteristics	101
Index	105

Using This Documentation

- **Overview** – Describes how to write applications that partition and manage system resources, such as processor sets and scheduling class, by referencing appropriate programming APIs.
- **Audience** – Developers writing resource management applications for the Oracle Solaris operating system.
- **Required knowledge** – Familiarity with C programming. Experience with virtualized environments is a plus.

Product Documentation Library

Documentation and resources for this product and related products are available at <http://www.oracle.com/pls/topic/lookup?ctx=E53394-01>.

Feedback

Provide feedback about this documentation at <http://www.oracle.com/goto/docfeedback>.

◆◆◆ 1 CHAPTER 1

Resource Management in the Oracle Solaris Operating System

Resource management can help developers who are writing either utility applications for managing computer resources or self-monitoring applications that can check their own usage and adjust accordingly. This chapter provides an introduction to resource management in the Oracle Solaris operating system. It covers the following topics:

- [“Understanding Resource Management in the Oracle Solaris Operating System” on page 11](#)
- [“Writing Resource Management Applications” on page 14](#)

Understanding Resource Management in the Oracle Solaris Operating System

The main concept behind resource management is that workloads on a server must be balanced for the system to work efficiently. Without good resource management, faulty runaway workloads can bring progress to a halt, causing unnecessary delays to priority jobs. Efficient resource management also enables organizations to economize by consolidating systems.

The Oracle Solaris operating system provides a structure for organizing workloads and resources, and provides controls for defining the quantity of resources that a particular unit of workload can consume. For an in-depth discussion of resource management from the system administrator's viewpoint, see [Chapter 1, “Introduction to Resource Management” in *Administering Resource Management in Oracle Solaris 11.3*](#).

Resource Management Workload Organization

The basic unit of workload is the *process*. Process IDs (PIDs) are numbered sequentially throughout the system. By default, each user is assigned by the system administrator to a

project, which is a network-wide administrative identifier. Each successful login to a project creates a new *task*, which is a grouping mechanism for processes. A task contains the login process as well as subsequent child processes.

For more information about projects and tasks, see [Chapter 2, “About Projects and Tasks” in *Administering Resource Management in Oracle Solaris 11.3*](#) for the system administrator's perspective or [Chapter 2, “Workload Hierarchy of Projects and Tasks”](#) in this document for the developer's point of view.

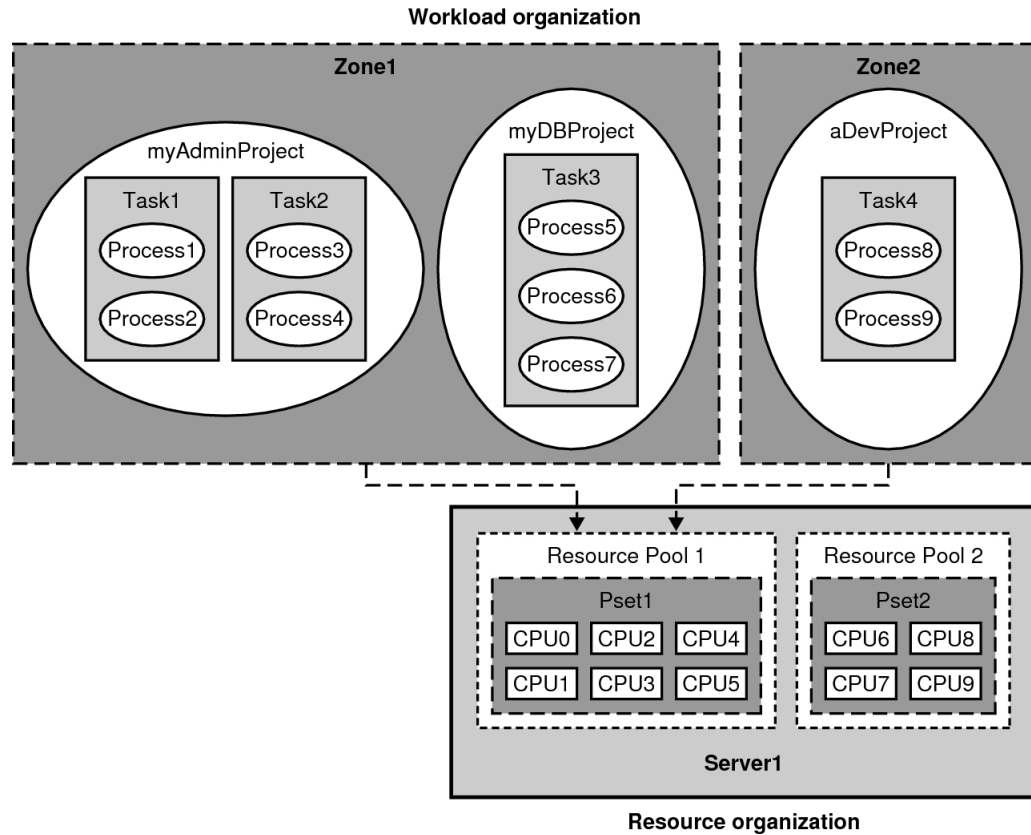
Processes can optionally be grouped into *non-global zones*, which are set up by system administrators for security purposes and to isolate processes. A zone enables one or more applications to run isolated from all other applications on the system. Non-global zones are discussed thoroughly in [Creating and Using Oracle Solaris Zones](#). To learn more about special precautions for writing resource management applications that run in zones, see [Chapter 7, “Design Considerations for Resource Management Applications in Oracle Solaris Zones”](#).

Resource Organization

The system administrator can assign workloads to specific CPUs or defined groups of CPUs in the system. CPUs can be grouped into *processor sets*, otherwise known as *psets*. A pset in turn can be coupled with one or more thread scheduling classes, which define CPU priorities, into a *resource pool*. Resource pools provide a convenient mechanism for a system administrator to make system resources available to users. [Chapter 12, “About Resource Pools” in *Administering Resource Management in Oracle Solaris 11.3*](#) covers resource pools for system administrators. Programming considerations are described in [Chapter 6, “Resource Pools in Oracle Solaris”](#).

The following diagram illustrates how workload and computer resources are organized in the Oracle Solaris operating system.

FIGURE 1 Workload and Resource Organization in the Oracle Solaris Operating System



Resource Controls

Simply assigning a workload unit to a resource unit is insufficient for managing the quantity of resources that users consume. To manage resources, the Oracle Solaris operating system provides a set of flags, actions, and signals that are referred to collectively as *resource controls*. Resource controls are stored in the `/etc/project` file or in a zone's configuration through the `zonecfg` command described in [zonecfg\(1M\)](#). The Fair Share Scheduler (FSS), for example, can allocate shares of CPU resources among workloads based on the specified importance factor for the workloads. With these resource controls, a system administrator can set privilege

levels and limit definitions for a specific zone, project, task, or process. To learn how a system administrator uses resource controls, see [Chapter 6, “About Resource Controls” in *Administering Resource Management in Oracle Solaris 11.3*](#). For programming considerations, see [Chapter 5, “Resource Controls in Oracle Solaris”](#).

Extended Accounting Facility

In addition to workload and resource organization, the Oracle Solaris operating system provides the *extended accounting facility* for monitoring and recording system resource usage. The extended accounting facility provides system administrators with a detailed set of resource consumption statistics on processes and tasks.

The facility is described in depth for system administrators in [Chapter 4, “About Extended Accounting” in *Administering Resource Management in Oracle Solaris 11.3*](#). The Oracle Solaris operating system provides developers with both a C interface and a Perl interface to the extended accounting facility. Refer to [Chapter 3, “Using the C Interface to Extended Accounting”](#) for the C interface or [Chapter 4, “Using the Perl Interface to Extended Accounting”](#) for the Perl interface.

Writing Resource Management Applications

This guide focuses on resource management from the developer's point of view and presents information for writing the following kinds of applications:

- Resource administration applications – Utilities to perform such tasks as allocating resources, creating partitions, and scheduling jobs.
- Resource monitoring applications – Applications that check system statistics through `kstats` to determine resource usage by systems, workloads, processes, and users.
- Resource accounting utilities – Applications that provide accounting information for analysis, billing, and capacity planning.
- Self-adjusting applications – Applications that can determine their use of resources and can adjust consumption as necessary.

Workload Hierarchy of Projects and Tasks

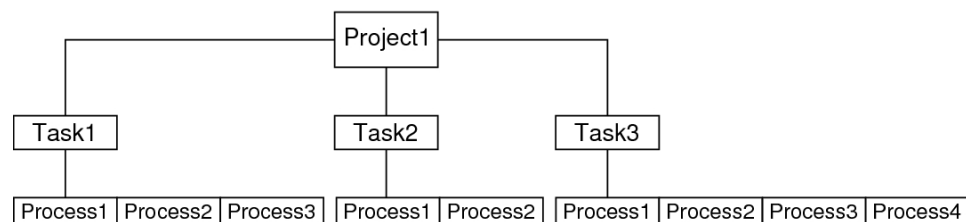
The chapter discusses the workload hierarchy and provides information about projects and tasks. It covers the following topics:

- [“Overview of Projects and Tasks” on page 15](#)
- [“Project and Task API Functions” on page 18](#)
- [“Code Examples for Accessing project Database Entries” on page 19](#)
- [“Programming Issues Associated With Projects and Tasks” on page 20](#)

Overview of Projects and Tasks

The Oracle Solaris operating system uses the workload hierarchy to organize the work being performed on the system. A *task* is a collection of processes that represents a workload component. A *project* is a collection of tasks that represents an entire workload. At any given time, a process can be a component of only one task and one project. The relationships in the workload hierarchy are illustrated in the following figure.

FIGURE 2 Workload Hierarchy



A user who is a member of more than one project can run processes in multiple projects at the same time.

All processes that are started by a process inherit the project and task created by the parent process. When you switch to a new project in a startup script, all child processes run in the new project.

An executing user process has an associated user identity (`uid`), group identity (`gid`), and project identity (`projid`). Process attributes and abilities are inherited from the user, group, and project identities to form the execution context for a task.

For an in-depth discussion of projects and tasks, see [Chapter 2, “About Projects and Tasks” in *Administering Resource Management in Oracle Solaris 11.3*](#). For the administration commands for managing projects and tasks, see [Chapter 3, “Administering Projects and Tasks” in *Administering Resource Management in Oracle Solaris 11.3*](#).

/etc/project File

The project file is the heart of workload hierarchy. The project database is maintained on a system through the `/etc/project` file or over the network through a naming service, such as NIS or LDAP.

The `/etc/project` file contains five standard projects.

`system`

This project is used for all system processes and daemons.

`user.root`

All root processes spawned by root logins and root cron, at, and batch jobs.

`noproject`

This special project is for IPQoS.

`default`

A default project is assigned to every user.

`group.staff`

This project is used for all users in the group `staff`.

To access the project file programmatically, use the following structure:

```
struct project {
    char    *pj_name;        /* name of the project */
```



```

projid_t  pj_projid;    /* numerical project ID */
char      *pj_comment; /* project comment */
char      **pj_users;  /* vector of pointers to project user names */
char      **pj_groups; /* vector of pointers to project group names */
char      *pj_attr;    /* project attributes */
};

```

The project structure members include the following:

`*pj_name`

Name of the project.

`projid`

Project ID.

`*pj_comment`

User-supplied project description.

`**pj_users`

Pointers to project user members.

`**pj_groups`

Pointers to project group members.

`*pj_attr`

Project attributes. Use these attributes to set values for resource controls and project pools.

Resource usage can be controlled through project attributes, or, for zones, configured through the `zonectg` command. Four prefixes are used to group the types of resource control attributes:

- `project.*` – This prefix denotes attributes that are used to control projects. For example, `project.max-locked-memory` indicates the total amount of locked memory allowed, expressed as a number of bytes. The `project.pool` attribute binds a project to a resource pool. See [Chapter 6, “Resource Pools in Oracle Solaris”](#).
- `task.*` – This prefix is used for attributes that are applied to tasks. For example, the `task.max-cpu-time` attribute sets the maximum CPU time that is available to this task's processes, expressed as a number of seconds.
- `process.*` – This prefix is used for process controls. For example, the `process.max-file-size` control sets the maximum file offset that is available for writing by this process, expressed as a number of bytes.
- `zone.*` – The `zone.*` prefix indicates a zone-wide resource control applied to projects, tasks, and processes in a zone. For example, `zone.max-lwps` prevents too many LWPs in

one zone from affecting other zones. A zone's total LWPs can be further subdivided among projects within the zone within the zone by using `project.max-lwps` entries.

For the complete list of resource controls, see [resource-controls\(5\)](#).

Project and Task API Functions

The following functions are provided to assist developers in working with projects. The functions use entries that describe user projects in the project database.

See *man pages section 3: Extended Library Functions, Volume 3* for additional information on these functions.

`endproject(3PROJECT)`

Close the project database and deallocate resources when processing is complete.

`fgetproject(3PROJECT)`

Returns a pointer to a structure containing an entry in the project database. Rather than using `nsswitch.conf`, `fgetproject()` reads a line from a stream.

`getdefaultproj(3PROJECT)`

Check the validity of the project keyword, look up the project, and return a pointer to the project structure if found.

`getprojbyid(3PROJECT)`

Search the project database for an entry with the number that specifies the project ID.

`getprojbyname(3PROJECT)`

Search the project database for an entry with the string that specifies project name.

`getproject(3PROJECT)`

Returns a pointer to a structure containing an entry in the project database.

`inproj(3PROJECT)`

Check whether the specified user is permitted to use the specified project.

`setproject(3PROJECT)`

Calling process joins the target project by creating a new task in the target project.

`setproject(3PROJECT)`

Rewind the project database to allow repeated searches.

Reentrant functions `getproject()`, `getprojbyname()`, `getprojbyid()`, `getdefaultproj()`, and `inproj()` use buffers supplied by the caller to store returned results. These functions are safe for use in both single-threaded applications and multithreaded applications.

Reentrant functions require three additional arguments:

- `proj`
- `buffer`
- `bufsize`

The `proj` argument must be a pointer to a `project` structure allocated by the caller. On successful completion, these functions return the project entry in this structure. Storage referenced by the `project` structure is allocated from the memory specified by the `buffer` argument. `bufsize` specifies the size in number of bytes.

If an incorrect buffer size is used, `getproject()` returns `NULL` and sets `errno` to `ERANGE`.

Code Examples for Accessing project Database Entries

EXAMPLE 1 Printing the First Three Fields of Each Entry in the project Database

Note the following key points for this example:

- `setproject()` rewinds the project database to start at the beginning.
- `getproject()` is called with a conservative maximum buffer size that is defined in `project.h`.
- `endproject()` closes the project database and frees resources.

```
#include <project.h>

struct project project;
char buffer[PROJECT_BUFSZ]; /* Use safe buffer size from project.h */

struct project *pp;

setproject(); /* Rewind the project database to start at the beginning */

while (1) {
    pp = getproject(&project, buffer, PROJECT_BUFSZ);
    if (pp == NULL)
        break;
    printf("%s:%d:%s\n", pp->pj_name, pp->pj_projid, pp->pj_comment);
}
```

```
};

endproject(); /* Close the database and free project resources */
```

EXAMPLE 2 Getting a project Database Entry That Matches the Caller's Project ID

The following example calls `getprojbyid()` to get a project database entry that matches the caller's project ID. The example then prints the project name and the project ID.

```
#include <project.h>

struct project *pj;
char buffer[PROJECT_BUFSZ]; /* Use safe buffer size from project.h */

main()
{
    projid_t pjid;
    pjid = getprojid();
    pj = getprojbyid(pjid, &project, buffer, PROJECT_BUFSZ);
    if (pj == NULL) {
        /* fail; */
    }
    printf("My project (name, id) is (%s, %d)\n", pp->pj_name, pp->pj_projid);
}
```

Programming Issues Associated With Projects and Tasks

Consider the following issues when writing your application:

- No function exists to explicitly create a new project.
- A user cannot log in if no default project for the user exists in the project database.
- A new task in the user's default project is created when the user logs in.
- When a process joins a project, the project's resource control and pool settings are applied to the process.
- `setproject()` requires privilege. The `newtask` command does not require privilege if you own the process. Either method can be used to create a task, but only `newtask` can change the project of a running process.
- No parent/child relationship exists between tasks.
- Finalized tasks can be created by using `newtask -F` or by using `setproject()` to associate the caller with a new project. Finalized tasks are useful when trying to accurately estimate aggregate resource accounting.

Using the C Interface to Extended Accounting

This chapter describes the C interface to extended accounting. It covers the following topics:

- [“Overview of the C Interface to Extended Accounting” on page 21](#)
- [“Extended Accounting API Functions” on page 21](#)
- [“C Code Examples for Accessing exacct Files” on page 24](#)

Overview of the C Interface to Extended Accounting

The extended accounting subsystem monitors resource consumption by workloads that are running on the system. Extended accounting produces accounting records for the workload tasks and processes.

For an overview of extended accounting and example procedures for administering extended accounting, see [Chapter 4, “About Extended Accounting” in *Administering Resource Management in Oracle Solaris 11.3*](#) and [Chapter 5, “Administering Extended Accounting Tasks” in *Administering Resource Management in Oracle Solaris 11.3*](#).

The extended accounting framework has been expanded for zones. Each zone has its own extended accounting files for task and process-based accounting. The extended accounting files in the global zone contain accounting records for the global zone and for all non-global zones. The accounting records contain a zone name tag that the global zone administrator can use during the extraction of per-zone accounting data from the accounting files in the global zone.

Extended Accounting API Functions

The extended accounting API contains functions that perform actions such as:

- `exacct` system calls
- Operations on the `exacct` file
- Operations on `exacct` objects

This section provides tables describing these functions. The function name is a link to its man page.

exacct System Calls

The following table lists the system calls that interact with the extended accounting subsystem.

TABLE 1 Extended Accounting System Calls

Function	Description
putacct(2)	Provides privileged processes with the ability to tag accounting records with additional data that is specific to the process
getacct(2)	Enables privileged processes to request extended accounting buffers from the kernel for currently executing tasks and processes
wracct(2)	Requests the kernel to write resource usage data for a specified task or process

Operations on the exacct File

The following table lists the functions that provide access to the exacct file.

TABLE 2 exacct File Functions

Function	Description
ea_open(3EXACCT)	Opens an exacct file.
ea_close(3EXACCT)	Closes an exacct file.
ea_get_object(3EXACCT)	First time use on a group of objects reads data into an <code>ea_object_t</code> structure. Subsequent use on the group cycles through the objects in the group.
ea_write_object(3EXACCT)	Appends the specified object to the open exacct file.
ea_next_object(3EXACCT)	Reads the basic fields (<code>eo_catalog</code> and <code>eo_type</code>) into an <code>ea_object_t</code> structure and rewinds to the head of the record.
ea_previous_object(3EXACCT)	Skips back one object in the exacct file and reads the basic fields (<code>eo_catalog</code> and <code>eo_type</code>) into an <code>ea_object_t</code> .
ea_get_hostname(3EXACCT)	Gets the name of the host on which the exacct file was created.
ea_get_creator(3EXACCT)	Determines the creator of the exacct file.

Operations on exacct Objects

The following table lists the functions that are used to access exacct objects.

TABLE 3 exacct Object Functions

Function	Description
ea_set_item(3EXACCT)	Assigns an exacct object and sets its values.
ea_set_group(3EXACCT)	Sets the values of a group of exacct objects.
ea_match_object_catalog(3EXACCT)	Checks an exacct object's mask to see whether that object has a specific catalog tag.
ea_attach_to_object(3EXACCT)	Attaches an exacct object to a specified exacct object.
ea_attach_to_group(3EXACCT)	Attaches a chain of exacct objects as member items of a specified group.
ea_free_item(3EXACCT)	Frees the value fields in the specified exacct object.
ea_free_object(3EXACCT)	Frees the specified exacct object and any attached hierarchies of objects.

Extended Accounting Memory Management

The following table lists the functions associated with extended accounting memory management.

TABLE 4 Extended Accounting Memory Management Functions

Link to man page	Description
ea_pack_object(3EXACCT)	Converts an exacct object from unpacked (in-memory) representation to packed (in-file) representation.
ea_unpack_object(3EXACCT)	Converts an exacct object from packed (in-file) representation to unpacked (in-memory) representation.
ea_strdup(3EXACCT)	Duplicates a string that is to be stored inside an <code>ea_object_t</code> structure.
ea_strfree(3EXACCT)	Frees a string previously copied by <code>ea_strdup()</code> .
ea_alloc(3EXACCT)	Allocates a block of memory of the requested size. This block can be safely passed to <code>libexacct</code> functions, and can be safely freed by any of the <code>ea_free</code> functions.
ea_free(3EXACCT)	Frees a block of memory previously allocated by <code>ea_alloc()</code> .
ea_free_object(3EXACCT)	Frees variable-length data in object hierarchy.
ea_free_item(3EXACCT)	Frees value fields of the designated object if <code>EUP_ALLOC</code> is specified. The object is not freed. <code>ea_free_object()</code> frees the specified object and any attached hierarchy of objects. If the flag argument is set to <code>EUP_ALLOC</code> , <code>ea_free_object()</code> also frees any variable-length data in the object hierarchy. If the flag argument is set to <code>EUP_NOALLOC</code> , <code>ea_free_object()</code> does not free the variable-length data. In particular, these flags should correspond to the flags specified in calls to <code>ea_unpack_object(3EXACCT)</code> .
ea_copy_object(3EXACCT)	Copies an <code>ea_object_t</code> . If the source object is part of a chain, only the current object is copied. If the source object is a group, only the group object is copied without its list of members. The group object

Link to man page	Description
ea_copy_object_tree(3EXACCT)	eg_nobjs and eg_objs fields are set to 0 and NULL respectively. Use ea_copy_tree() to copy recursively a group or a list of items. ea_copy_object_tree recursively copies an ea_object_t. All elements in the eo_next list are copied. Any group objects are recursively copied. The returned object can be completely freed with ea_free_object(3EXACCT) by specifying the EUP_ALLOC flag.
ea_get_object_tree()	Reads in nobj top-level objects from the file, returning the same data structure that would have originally been passed to ea_write_object(). On encountering a group object, ea_get_object() reads only the group header part of the group. ea_get_object_tree() reads the group and all its member items, recursing into subrecords if necessary. The returned object data structure can be completely freed with ea_free_object() by specifying the EUP_ALLOC flag.

Extended Accounting Miscellaneous Operations

These functions are associated with miscellaneous operations:

- [ea_error\(3EXACCT\)](#) – Returns the error value of the last failure recorded by the invocation of one of the functions of the extended accounting library, libexacct.
- [ea_match_object_catalog\(3EXACCT\)](#) – Returns TRUE if the exacct object specified by obj has a catalog tag that matches the mask specified by catmask.

C Code Examples for Accessing exacct Files

This section provides code examples for accessing exacct files.

EXAMPLE 3 Displaying exacct Data for a Designated pid

This example displays a specific pid's exacct data snapshot from the kernel.

```
...
ea_object_t *scratch;
int unpk_flag = EUP_ALLOC; /* use the same allocation flag */
                          /* for unpack and free */

/* Omit return value checking, to keep code samples short */

bsize = getacct(P_PID, pid, NULL, 0);
```



```

buf = malloc(bsize);

/* Retrieve exacct object and unpack */
getacct(P_PID, pid, buf, bsize);
ea_unpack_object(&scratch, unpk_flag, buf, bsize);

/* Display the exacct record */
disp_obj(scratch);
if (scratch->eo_type == EO_GROUP) {
    disp_group(scratch);
}
ea_free_object(scratch, unpk_flag);
...

```

EXAMPLE 4 Identifying Individual Tasks During a Kernel Build

This example evaluates kernel builds and displays a string that describes the portion of the source tree being built by this task. It displays the portion of the source being built to aid in the per-source-directory analysis.

Note the following key points for this example:

- To aggregate the time for a make, which could include many processes, each make is initiated as a task. Child make processes are created as different tasks. To aggregate across the makefile tree, the parent-child task relationship must be identified.
- Add a tag with this information to the task's exacct file. Add a current working directory string that describes the portion of the source tree being built by this task's make operation.

```

ea_set_item(&cwd, EXT_STRING | EXC_LOCAL | MY_CWD,
           cwdbuf, strlen(cwdbuf));

...
/* Omit return value checking and error processing */
/* to keep code sample short */
ptid = gettaskid(); /* Save "parent" task-id */
tid = settaskid(getprojid(), TASK_NORMAL); /* Create new task */

/* Set data for item objects ptskid and cwd */
ea_set_item(&ptskid, EXT_UINT32 | EXC_LOCAL | MY_PTID, &ptid, 0);
ea_set_item(&cwd, EXT_STRING | EXC_LOCAL | MY_CWD, cwdbuf, strlen(cwdbuf));

/* Set grp object and attach ptskid and cwd to grp */
ea_set_group(&grp, EXT_GROUP | EXC_LOCAL | EXD_GROUP_HEADER);
ea_attach_to_group(&grp, &ptskid);
ea_attach_to_group(&grp, &cwd);

/* Pack the object and put it back into the accounting stream */

```

```
ea_buflen = ea_pack_object(&grp, ea_buf, sizeof(ea_buf));
putacct(P_TASKID, tid, ea_buf, ea_buflen, EP_EXACCT_OBJECT);

/* Memory management: free memory allocate in ea_set_item */
ea_free_item(&cwd, EUP_ALLOC);
...
```

EXAMPLE 5 Reading and Displaying the Contents of a System exacct File

This example shows how to read and display a system exacct file for a process or a task.

Note the following key points for this example:

- Calls `ea_get_object()` to get the next object in the file. Calls `ea_get_object()` in a loop until EOF enables a complete traversal of the exacct file.
- `catalog_name()` uses the `catalog_item` structure to convert an Oracle Solaris catalog's type ID to a meaningful string that describes the content of the object's data. The type ID is obtained by masking the lowest 24 bits, or 3 bytes.

```
switch(o->eo_catalog & EXT_TYPE_MASK) {
    case EXT_UINT8:
        printf(" 8: %u", o->eo_item.ei_uint8);
        break;
    case EXT_UINT16:
        ...
}
```

- The upper 4 bits of `TYPE_MASK` are used to find out the data type to print the object's actual data.
- `disp_group()` takes a pointer to a group object and the number of objects in the group. For each object in the group, `disp_group()` calls `disp_obj()` and recursively calls `disp_group()` if the object is a group object.

```
/* Omit return value checking and error processing */
/* to keep code sample short */
main(int argc, char *argv)
{
    ea_file_t ef;
    ea_object_t scratch;
    char *fname;

    fname = argv[1];
    ea_open(&ef, fname, NULL, EO_NO_VALID_HDR, O_RDONLY, 0);
    bzero(&scratch, sizeof(ea_object_t));
    while (ea_get_object(&ef, &scratch) != -1) {
        disp_obj(&scratch);
    }
}
```

```

        if (scratch.eo_type == EO_GROUP)
            disp_group(&ef, scratch.eo_group.eg_nobjs);
        bzero(&scratch, sizeof (ea_object_t));
    }
    ea_close(&ef);
}

struct catalog_item { /* convert Oracle Solaris catalog's type ID */
                    /* to a meaningful string */

    int type;
    char *name;
} catalog[] = {
    { EXD_VERSION, "version\t" },
    ...
    { EXD_PROC_PID, " pid\t" },
    ...
};

static char *
catalog_name(int type)
{
    int i = 0;
    while (catalog[i].type != EXD_NONE) {
        if (catalog[i].type == type)
            return (catalog[i].name);
        else
            i++;
    }
    return ("unknown\t");
}

static void disp_obj(ea_object_t *o)
{
    printf("%s\t", catalog_name(o->eo_catalog & 0xffffffff));
    switch(o->eo_catalog & EXT_TYPE_MASK) {
    case EXT_UINT8:
        printf(" 8: %u", o->eo_item.ei_uint8);
        break;
    case EXT_UINT16:
        ...
    }
}

static void disp_group(ea_file_t *ef, uint_t nobjs)
{
    for (i = 0; i < nobjs; i++) {
        ea_get_object(ef, &scratch);
        disp_obj(&scratch);
        if (scratch.eo_type == EO_GROUP)
            disp_group(ef, scratch.eo_group.eg_nobjs);
    }
}

```

```
    }  
}
```

Programming Issues With `exacct` Files

- Note the following issues related to memory management:
 - Use the same allocation flags for `ea_free_object()` and `ea_unpack_object()`.
 - For string objects, an `ea_set_item()` results in allocation, and should be followed by `ea_free_item(obj, EUP_ALLOC)` to free internal storage.
 - `ea_pack_object()` and `getacct()` use zero size. To get size, `getacct()` should be called twice: first time with `NULL` buffer to size buffer to be passed in the second call. See Example 3-1 in “[C Code Examples for Accessing `exacct` Files](#)” on page 24.
- In order to be robust in the face of changes to `exacct` file content, applications should skip unknown `exacct` records in `exacct` files produced by the system.
- Use `EXC_LOCAL` for customized accounting to create application-specific records. Use `libexacct` as general tracing or debugging facility.
 - See `<sys/exacct_catalog.h>`.
 - You can customize the data id field of `ea_catalog_t`.

Using the Perl Interface to Extended Accounting

The Perl interface provides a Perl binding to the extended accounting tasks and projects. The interface allows the accounting files produced by the `exacct` framework to be read by Perl scripts. The interface also allows the writing of `exacct` files by Perl scripts.

This chapter includes the following topics:

- [“Extended Accounting Overview” on page 29](#)
- [“Perl Interface to `libexacct`” on page 30](#)
- [“Perl Modules” on page 31](#)
- [“Perl Code Examples” on page 43](#)
- [“Output From the Perl `dump` Method” on page 47](#)

Extended Accounting Overview

Extended accounting (`exacct`) is an accounting framework for the Oracle Solaris operating system that provides additional functionality to that provided by the traditional SVR4 accounting mechanism. Traditional SVR4 accounting has these drawbacks:

- The data collected by traditional accounting cannot be modified.
The type or quantity of statistics SVR4 accounting gathers cannot be customized for each application. Changes to the data that traditional accounting collects would not work with all of the existing applications that use the accounting files.
- The SVR4 accounting mechanism is not open.
Applications cannot embed their own data in the system accounting data stream.
- The traditional accounting mechanism has no aggregation facilities.
The Oracle Solaris operating system writes an individual record for each process that exists. No facilities are provided for grouping sets of accounting records into higher-level aggregates.

The exacct framework addresses the limitations of traditional accounting and provides a configurable, open, and extensible framework for the collection of accounting data.

- The data that is collected can be configured using the exacct API.
- Applications can either embed their own data inside the system accounting files, or create and manipulate their own custom accounting files.
- The lack of data aggregation facilities in the traditional accounting mechanism are addressed by *tasks* and *projects*. Tasks identify a set of processes that are to be considered as a unit of work. Projects allow the processes executed by a set of users to be aggregated into a higher-level entity. See the [project\(4\)](#) man page for more details about tasks and projects.

For a more extensive overview of extended accounting, see [Chapter 4, “About Extended Accounting” in *Administering Resource Management in Oracle Solaris 11.3*](#).

Perl Interface to libexacct

Perl Interface Object Model

The `Sun::Solaris::Exacct` module is the parent of all the classes provided by the `libexacct(3LIB)` library. `libexacct(3LIB)` provides operations on types of entities: exacct format files, catalog tags, and exacct objects. exacct objects are subdivided into two types.

- Items – Single data values
- Groups – Lists of items

Benefits of Using the Perl Interface to libexacct

The Perl extensions to extended accounting provide a Perl interface to the underlying `libexacct(3LIB)` API and offer the following enhancements:

- Full equivalence to the C API that provides a Perl interface that is functionally equivalent to the underlying C API.

The interface provides a mechanism for accessing exacct files that does not require C coding. All the functionality that is available from C is also available by using the Perl interface.

- Ease of use.

Data obtained from the underlying C API is presented as Perl data types. Perl data types ease access to the data and remove the need for buffer pack and unpack operations.

- Automated memory management.

The C API requires the programmer to take responsibility for managing memory when accessing `exacct` files. Memory management involves passing the appropriate flags to functions, such as `ea_unpack_object(3EXACCT)`, and explicitly allocating buffers to pass to the API. The Perl API removes these requirements because all memory management is performed by the Perl library.

- Prevent incorrect use of API.

The `ea_object_t` structure provides the in-memory representation of `exacct` records. The `ea_object_t` structure is a union type that is used for manipulating both Group and Item records. As a result, an incorrectly typed structure can be passed to some of the API functions. The addition of a class hierarchy prevents this type of programming error.

Perl Double-Typed Scalars

The modules described in this document make extensive use of the Perl double-typed scalar facility. The *double-typed scalar* facility allows a scalar value to behave either as an integer or as a string depending upon the context. This behavior is the same as exhibited by the `#!` Perl variable (`errno`). The double-typed scalar facility avoids the need to map from an integer value into the corresponding string in order to display a value. The following example illustrates the use of double-typed scalars.

```
# Assume $obj is a Sun::Solaris::Item
my $type = $obj->type();

# prints out "2 EO_ITEM"
printf("%d %s\n", $type, $type);

# Behaves as an integer, $i == 2
my $i = 0 + $type;

# Behaves as a string, $s = "abc EO_ITEM xyx"
my $s = "abc $type xyz";
```

Perl Modules

The Perl modules each contain a group of related project, task and `exacct`-related functions. Each function has the standard `Sun::Solaris::` Perl package prefix.

TABLE 5 Perl Modules

Module	Description
“Sun::Solaris::Project Module” on page 33	Provides functions to access the project manipulation functions: <code>getprojid(2)</code> , <code>setproject(3PROJECT)</code> , <code>project_walk(3PROJECT)</code> , <code>getproject(3PROJECT)</code> , <code>getprojbyname(3PROJECT)</code> , <code>getprojbyid(3PROJECT)</code> , <code>getdefaultproj(3PROJECT)</code> , <code>inproj(3PROJECT)</code> , <code>getprojidbyname(3PROJECT)</code> , <code>setproject(3PROJECT)</code> , <code>endproject(3PROJECT)</code> , <code>fgetproject(3PROJECT)</code> .
“Sun::Solaris::Task Module” on page 34	Provides functions to access the task manipulation functions <code>settaskid(2)</code> and <code>gettaskid(2)</code> .
“Sun::Solaris::Exacct Module” on page 35	Top-level exacct module. Functions in this module access both the exacct-related system calls <code>getacct(2)</code> , <code>putacct(2)</code> , and <code>wracct(2)</code> as well as the <code>libexacct(3LIB)</code> library function <code>ea_error(3EXACCT)</code> . This module contains constants for all the various exacct <code>E0_*</code> , <code>EW_*</code> , <code>EXR_*</code> , <code>P_*</code> and <code>TASK_*</code> macros.
“Sun::Solaris::Exacct::Catalog Module” on page 36	Provides object-oriented methods to access the bitfields within an exacct catalog tag as well as the <code>EXC_*</code> , <code>EXD_*</code> and <code>EXD_*</code> macros.
“Sun::Solaris::Exacct::File Module” on page 38	Provides object-oriented methods to access the <code>libexacct(3LIB)</code> accounting file functions: <code>ea_open(3EXACCT)</code> , <code>ea_close(3EXACCT)</code> , <code>ea_get_creator(3EXACCT)</code> , <code>ea_get_hostname(3EXACCT)</code> , <code>ea_next_object(3EXACCT)</code> , <code>ea_previous_object(3EXACCT)</code> , <code>ea_write_object(3EXACCT)</code> .
“Sun::Solaris::Exacct::Object Module” on page 40	Provides object-oriented methods to access the individual exacct accounting file object. An exacct object is represented as an opaque reference that is blessed into the appropriate <code>Sun::Solaris::Exacct::Object</code> subclass. This module is further subdivided into the two types of possible object: <code>Item</code> and <code>Group</code> . Methods are also provided to access the <code>ea_match_object_catalog(3EXACCT)</code> , <code>ea_attach_to_object(3EXACCT)</code> functions.
“Sun::Solaris::Exacct::Object::Item Module” on page 41	Provides object-oriented methods to access an individual exacct accounting file <code>Item</code> . Objects of this type inherit from <code>Sun::Solaris::Exacct::Object</code> .

Module	Description
“Sun::Solaris::Exacct::Object::Group Module” on page 42	Provides object-oriented methods to access an individual exacct accounting file Group. Objects of this type inherit from Sun::Solaris::Exacct::Object, and provide access to the ea_attach_to_group(3EXACCT) function. The Items contained within the Group are presented as a perl array.
“Sun::Solaris::Exacct::Object::_Array Module” on page 43	Private array type, used as the type of the array within a Sun::Solaris::Exacct::Object::Group.

Sun::Solaris::Project Module

The Sun::Solaris::Project module provides wrappers for the project-related system calls and the libproject(3LIB) library.

Sun::Solaris::Project Constants

The Sun::Solaris::Project module uses constants from the project-related header files.

```
MAXPROJID
PROJNAME_MAX
PROJF_PATH
PROJECT_BUFSZ
SETPROJ_ERR_TASK
SETPROJ_ERR_POOL
```

Sun::Solaris::Project Functions, Class Methods, and Object Methods

The perl extensions to the libexacct(3LIB) API provide the following functions for projects:

```
setproject(3PROJECT)           fgetproject(3PROJECT)
setproject(3PROJECT)           getprojbyname(3PROJECT)
getdefaultproj(3PROJECT)       getprojbyid(3PROJECT)
inproj(3PROJECT)               getprojbyname(3PROJECT)
getproject(3PROJECT)           endproject(3PROJECT)
```

The Sun::Solaris::Project module has no class methods.

The Sun::Solaris::Project module has no object methods.

Sun::Solaris::Project Exports

By default, nothing is exported from this module. You can use the tags listed in the following table to selectively import the constants and functions defined in this module.

Tag	Constant or Function
:SYSCALLS	getprojid()
:LIBCALLS	setproject(), activeprojects(), getprojent(), setprojent(), endprojent(), getprojbyname(), getprojbyid(), getdefaultproj(), fgetprojent(), inproj(), getprojidbyname()
:CONSTANTS	MAXPROJID_TASK, PROJNAME_MAX, PROJF_PATH, PROJECT_BUFSZ, SETPROJ_ERR, SETPROJ_ERR_POOL
:ALL	:SYSCALLS, :LIBCALLS, :CONSTANTS

Sun::Solaris::Task Module

The Sun::Solaris::Task module provides wrappers for the `settaskid(2)` and `gettaskid(2)` system calls.

Sun::Solaris::Task Constants

The Sun::Solaris::Task module uses the following constants:

TASK_NORMAL
TASK_FINAL

Sun::Solaris::Task Functions, Class Methods, and Object Methods

The perl extensions to the `libexecct(3LIB)` API provides the following functions for tasks:

[settaskid\(2\)](#)
[gettaskid\(2\)](#)

The Sun::Solaris::Task module has no class methods.

The Sun::Solaris::Task module has no object methods.

Sun::Solaris::Task Exports

By default, nothing is exported from this module. You can use the tags listed in the following table to selectively import the constants and functions defined in this module.

Tag	Constant or Function
:SYSCALLS	settaskid(), gettaskid()
:CONSTANTS	TASK_NORMAL and TASK_FINAL
:ALL	:SYSCALLS and :CONSTANTS

Sun::Solaris::Exacct Module

The Sun::Solaris::Exacct module provides wrappers for the `ea_error(3EXACCT)` function and for all the exacct system calls.

Sun::Solaris::Exacct Constants

The Sun::Solaris::Exacct module provides constants from the various exacct header files. The `P_PID`, `P_TASKID`, `P_PROJID` and all the `EW_*`, `EP_*`, `EXR_*` macros are extracted during the module build process. The macros are extracted from the exacct header files under `/usr/include` and provided as Perl constants. Constants passed to the Sun::Solaris::Exacct functions can either be an integer value such as `EW_FINAL` or a string representation of the same variable such as `"EW_FINAL"`.

Sun::Solaris::Exacct Functions, Class Methods, and Object Methods

The perl extensions to the `libexacct(3LIB)` API provide the following functions for the Sun::Solaris::Exacct module.

<code>getacct(2)</code>	<code>ea_register_catalog</code>
<code>putacct(2)</code>	<code>ea_new_file</code>
<code>wracct(2)</code>	<code>ea_new_item</code>
<code>ea_error(3EXACCT)</code>	<code>ea_new_group</code>
<code>ea_error_str</code>	<code>ea_dump_object</code>

Note - `ea_error_str()` is provided as a convenience, so that repeated blocks of code like the following can be avoided:

```
if (ea_error() == EXR_SYSCALL_FAIL) {
    print("error: $!\n");
} else {
    print("error: ", ea_error(), "\n");
}
```

The `Sun::Solaris::Exacct` module has no class methods.

The `Sun::Solaris::Exacct` module has no object methods.

Sun::Solaris::Exacct Exports

By default, nothing is exported from this module. The following tags can be used to selectively import constants and functions defined in this module.

Tag	Constant or Function
:SYSCALLS	<code>getacct()</code> , <code>putacct()</code> , <code>wracct()</code>
:LIBCALLS	<code>ea_error()</code> , <code>ea_error_str()</code>
:CONSTANTS	<code>P_PID</code> , <code>P_TASKID</code> , <code>P_PROJID</code> <code>EW_*</code> , <code>EP_*</code> , <code>EXR_*</code>
:SHORTAND	<code>ea_register_catalog()</code> , <code>ea_new_catalog()</code> , <code>ea_new_file()</code> , <code>ea_new_item()</code> , <code>ea_new_group()</code> , <code>ea_dump_object()</code>
:ALL	:SYSCALLS, :LIBCALLS, :CONSTANTS and :SHORTAND
:EXACCT_CONSTANTS	:CONSTANTS, plus the :CONSTANTS tags for <code>Sun::Solaris::Catalog</code> , <code>Sun::Solaris::File</code> , <code>Sun::Solaris::Object</code>
:EXACCT_ALL	:ALL, plus the :ALL tags for <code>Sun::Solaris::Catalog</code> , <code>Sun::Solaris::File</code> , <code>Sun::Solaris::Object</code>

Sun::Solaris::Exacct::Catalog Module

The `Sun::Solaris::Exacct::Catalog` module provides a wrapper around the 32-bit integer used as a catalog tag. The catalog tag is represented as a Perl object blessed into the `Sun::Solaris::Exacct::Catalog` class. Methods can be used to manipulate fields in a catalog tag.

Sun::Solaris::Exacct::Catalog Constants

All the `EXT_*`, `EXC_*` and `EXD_*` macros are extracted during the module build process from the `/usr/include/sys/exact_catalog.h` file and are provided as constants. Constants passed to the `Sun::Solaris::Exacct::Catalog` methods can either be an integer value, such as `EXT_UINT8`, or the string representation of the same variable, such as `"EXT_UINT8"`.

Sun::Solaris::Exacct::Catalog Functions, Class Methods, and Object Methods

The Perl extensions to the `libexacct(3LIB)` API provide the following class methods for `Sun::Solaris::Exacct::Catalog`. [Exacct\(3PERL\)](#) and [Exacct::Catalog\(3PERL\)](#).

```
register
new
```

The Perl extensions to the `libexacct(3LIB)` API provide the following object methods for `Sun::Solaris::Exacct::Catalog`.

```
value                type_str
type                 catalog_str
catalog              id_str
id
```

Sun::Solaris::Exacct::Catalog Exports

By default, nothing is exported from this module. The following tags can be used to selectively import constants and functions defined in this module.

Tag	Constant or Function
:CONSTANTS	EXT_*, EXC_* and EXD_*.
:ALL	:CONSTANTS

Additionally, any constants defined with the `register()` function can optionally be exported into the caller's package.

Sun::Solaris::Exacct::File Module

The Sun::Solaris::Exacct::File module provides wrappers for the exacct functions that manipulate accounting files. The interface is object-oriented, and allows the creation and reading of exacct files. The C library calls that are wrapped by this module are:

```
ea_open(3EXACCT)           ea_write_object(3EXACCT)
ea_close(3EXACCT)         ea_get_object(3EXACCT)
ea_next_object(3EXACCT)   ea_get_creator(3EXACCT)
ea_previous_object(3EXACCT) ea_get_hostname(3EXACCT)
```

The file read and write methods operate on Sun::Solaris::Exacct::Object objects. These methods perform all the necessary memory management, packing, unpacking and structure conversions that are required.

Sun::Solaris::Exacct::File Constants

Sun::Solaris::Exacct::File provides the EO_HEAD, EO_TAIL, EO_NO_VALID_HDR, EO_POSN_MSK and EO_VALIDATE_MSK constants. Other constants that are needed by the new() method are in the standard Perl Fcntl module. [Table 6, “\\$oflags and \\$aflags Parameters,” on page 38](#) describes the action of new() for various values of \$oflags and \$aflags.

Sun::Solaris::Exacct::File Functions, Class Methods, and Object Methods

The Sun::Solaris::Exacct::File module has no functions.

The Perl extensions to the libexacct(3LIB) API provide the following class method for Sun::Solaris::Exacct::File.

new

The following table describes the new() action for combinations of the \$oflags and \$aflags parameters.

TABLE 6 \$oflags and \$aflags Parameters

\$oflags	\$aflags	Action
O_RDONLY	Absent or EO_HEAD	Open for reading at the start of the file.

\$oflags	\$aflags	Action
O_RDONLY	EO_TAIL	Open for reading at the end of the file.
O_WRONLY	Ignored	File must exist, open for writing at the end of the file.
O_WRONLY O_CREAT	Ignored	Create file if the file does not exist. Otherwise, truncate, and open for writing.
O_RDWR	Ignored	File must exist, open for reading or writing, at the end of the file.
O_RDWR O_CREAT	Ignored	Create file if the file does not exist. Otherwise, truncate, and open for reading or writing.

Note - The only valid values for `$oflags` are the combinations of `O_RDONLY`, `O_WRONLY`, `O_RDWR` or `O_CREAT`. `$aflags` describes the required positioning in the file for `O_RDONLY`. Either `EO_HEAD` or `EO_TAIL` are allowed. If absent, `EO_HEAD` is assumed.

The perl extensions to the `libexact(3LIB)` API provide the following object methods for `Sun::Solaris::Exact::File`.

```
creator
hostname
next
previous
get
write
```

Note - Close a `Sun::Solaris::Exact::File`. There is no explicit `close()` method for a `Sun::Solaris::Exact::File`. The file is closed when the `filehandle` object is undefined or reassigned.

Sun::Solaris::Exact::File Exports

By default, nothing is exported from this module. The following tags can be used to selectively import constants that are defined in this module.

Tag	Constant or Function
:CONSTANTS	EO_HEAD, EO_TAIL, EO_NO_VALID_HDR, EO_POSN_MSK, EO_VALIDATE_MSK.
:ALL	:CONSTANTS and <code>Fcntl(:DEFAULT)</code> .

Sun::Solaris::Exacct::Object Module

The Sun::Solaris::Exacct::Object module serves as a parent of the two possible types of exacct objects: Items and Groups. An exacct Item is a single data value, an embedded exacct object, or a block of raw data. An example of a single data value is the number of seconds of user CPU time consumed by a process. An exacct Group is an ordered collection of exacct Items such as all of the resource usage values for a particular process or task. If Groups need to be nested within each other, the inner Groups can be stored as embedded exacct objects inside the enclosing Group.

The Sun::Solaris::Exacct::Object module contains methods that are common to both exacct Items and Groups. Note that the attributes of Sun::Solaris::Exacct::Object and all classes derived from it are read-only after initial creation via `new()`. The attributes made read-only prevents the inadvertent modification of the attributes which could give rise to inconsistent catalog tags and data values. The only exception to the read-only attributes is the array used to store the Items inside a Group object. This array can be modified using the normal perl array operators.

Sun::Solaris::Exacct::Object Constants

Sun::Solaris::Exacct::Object provides the `EO_ERROR`, `EO_NONE`, `EO_ITEM` and `EO_GROUP` constants.

Sun::Solaris::Exacct::Object Functions, Class Methods, and Object Methods

The Sun::Solaris::Exacct::Object module has no functions.

The Perl extensions to the `libexacct(3LIB)` API provide the following class method for Sun::Solaris::Exacct::Object.

`dump`

The Perl extensions to the `libexacct(3LIB)` API provide the following object methods for Sun::Solaris::Exacct::Object.

`type`
`catalog`
`match_catalog`
`value`

Sun::Solaris::Exacct::Object Exports

By default, nothing is exported from this module. The following tags can be used to selectively import constants and functions defined in this module.

Tag	Constant or Function
:CONSTANTS	EO_ERROR, EO_NONE, EO_ITEM and EO_GROUP
:ALL	:CONSTANTS

Sun::Solaris::Exacct::Object::Item Module

The Sun::Solaris::Exacct::Object::Item module is used for exacct data Items. An exacct data Item is represented as an opaque reference, blessed into the Sun::Solaris::Exacct::Object::Item class, which is a subclass of the Sun::Solaris::Exacct::Object class. The underlying exacct data types are mapped onto Perl types as follows.

TABLE 7 exacct Data Types Mapped to Perl Data Types

exacct type	Perl internal type
EXT_UINT8	IV (integer)
EXT_UINT16	IV (integer)
EXT_UINT32	IV (integer)
EXT_UINT64	IV (integer)
EXT_DOUBLE	NV (double)
EXT_STRING	PV (string)
EXT_EXACCT_OBJECT	Sun::Solaris::Exacct::Object subclass
EXT_RAW	PV (string)

Sun::Solaris::Exacct::Object::Item Constants

Sun::Solaris::Exacct::Object::Item has no constants.

Sun::Solaris::Exacct::Object::Item Functions, Class Methods, and Object Methods

Sun::Solaris::Exacct::Object::Item has no functions.

`Sun::Solaris::Exacct::Object::Item` inherits all class methods from the `Sun::Solaris::Exacct::Object` base class, plus the `new()` class method.

`new`

`Sun::Solaris::Exacct::Object::Item` inherits all object methods from the `Sun::Solaris::Exacct::Object` base class.

Sun::Solaris::Exacct::Object::Item Exports

`Sun::Solaris::Exacct::Object::Item` has no exports.

Sun::Solaris::Exacct::Object::Group Module

The `Sun::Solaris::Exacct::Object::Group` module is used for `exacct Group` objects. An `exacct Group` object is represented as an opaque reference, blessed into the `Sun::Solaris::Exacct::Object::Group` class, which is a subclass of the `Sun::Solaris::Exacct::Object` class. The `Items` within a `Group` are stored inside a Perl array, and a reference to the array can be accessed via the inherited `value()` method. This means that the individual `Items` within a `Group` can be manipulated with the normal Perl array syntax and operators. All data elements of the array must be derived from the `Sun::Solaris::Exacct::Object` class. `Group` objects can also be nested inside each other merely by adding an existing `Group` as a data `Item`.

Sun::Solaris::Exacct::Object::Group Constants

`Sun::Solaris::Exacct::Object::Group` has no constants.

Sun::Solaris::Exacct::Object::Group Functions, Class Methods, and Object Methods

`Sun::Solaris::Exacct::Object::Group` has no functions.

`Sun::Solaris::Exacct::Object::Group` inherits all class methods from the `Sun::Solaris::Exacct::Object` base class, plus the `new()` class method.

`new`

`Sun::Solaris::Exacct::Object::Group` inherits all object methods from the `Sun::Solaris::Exacct::Object` base class, plus the `new()` class method.

as_hash
as_hashlist

Sun::Solaris::Exacct::Object::Group Exports

Sun::Solaris::Exacct::Object::Group has no exports.

Sun::Solaris::Exacct::Object::_Array Module

The Sun::Solaris::Exacct::Object::_Array class is used internally for enforcing type checking of the data Items that are placed in an exacct Group. Sun::Solaris::Exacct::Object::_Array should not be created directly by the user.

Sun::Solaris::Exacct::Object::_Array Constants

Sun::Solaris::Exacct::Object::_Array has no constants.

Sun::Solaris::Exacct::Object::_Array Functions, Class Methods, and Object Methods

Sun::Solaris::Exacct::Object::_Array has no functions.

Sun::Solaris::Exacct::Object::_Array has internal-use class methods.

Sun::Solaris::Exacct::Object::_Array uses perl TIEARRAY methods.

Sun::Solaris::Exacct::Object::_Array Exports

Sun::Solaris::Exacct::Object::_Array has no exports.

Perl Code Examples

This section shows perl code examples for accessing exacct files.

EXAMPLE 6 Using the Pseudocode Prototype

In typical use the Perl `exacct` library reads existing `exacct` files. Use pseudocode to show the relationships of the various Perl `exacct` classes. Illustrate in pseudocode the process of opening and scanning an `exacct` file, and processing objects of interest. In the following pseudocode, the ‘convenience’ functions are used in the interest of clarity.

```
-- Open the exacct file ($f is a Sun::Solaris::Exacct::File)
my $f = ea_new_file(...)

-- While not EOF ($o is a Sun::Solaris::Exacct::Object)
while (my $o = $f->get())

    -- Check to see if object is of interest
    if ($o->type() == &EO_ITEM)
        ...

    -- Retrieve the catalog ($c is a Sun::Solaris::Exacct::Catalog)
    $c = $o->catalog()

    -- Retrieve the value
    $v = $o->value();

    -- $v is a reference to a Sun::Solaris::Exacct::Group for a Group
    if (ref($v))
        ....

    -- $v is perl scalar for Items
    else
```

EXAMPLE 7 Recursively dumping an `exacct` Object

```
sub dump_object
{
    my ($obj, $indent) = @_;
    my $istr = ' ' x $indent;

    #
    # Retrieve the catalog tag. Because we are doing this in an array
    # context, the catalog tag will be returned as a (type, catalog, id)
    # triplet, where each member of the triplet will behave as an integer
    # or a string, depending on context. If instead this next line provided
    # a scalar context, e.g.
    #   my $cat = $obj->catalog()->value();
    # then $cat would be set to the integer value of the catalog tag.
    #
```

```
my @cat = $obj->catalog()->value();

#
# If the object is a plain item
#
if ($obj->type() == &EO_ITEM) {
    #
    # Note: The '%s' formats provide s string context, so the
    # components of the catalog tag will be displayed as the
    # symbolic values. If we changed the '%s' formats to '%d',
    # the numeric value of the components would be displayed.
    #
    printf("%sITEM\n%s Catalog = %s|%s|%s\n",
           $istr, $istr, @cat);
    $indent++;

    #
    # Retrieve the value of the item. If the item contains in
    # turn a nested exact object (i.e. a item or group), then
    # the value method will return a reference to the appropriate
    # sort of perl object (Exact::Object::Item or
    # Exact::Object::Group). We could of course figure out that
    # the item contained a nested item or group by examining
    # the catalog tag in @cat and looking for a type of
    # EXT_EXACCT_OBJECT or EXT_GROUP.
    my $val = $obj->value();
    if (ref($val)) {
        # If it is a nested object, recurse to dump it.
        dump_object($val, $indent);
    } else {
        # Otherwise it is just a 'plain' value, so display it.
        printf("%s Value = %s\n", $istr, $val);
    }

    #
    # Otherwise we know we are dealing with a group. Groups represent
    # contents as a perl list or array (depending on context), so we
    # can process the contents of the group with a 'foreach' loop, which
    # provides a list context. In a list context the value method
    # returns the content of the group as a perl list, which is the
    # quickest mechanism, but doesn't allow the group to be modified.
    # If we wanted to modify the contents of the group we could do so
    # like this:
    #   my $grp = $obj->value(); # Returns an array reference
    #   $grp->[0] = $newitem;
    # but accessing the group elements this way is much slower.
    #
    } else {
```

```
        printf("%sGROUP\n%s  Catalog = %s|%s|%s\n",
              $istr, $istr, @cat);
        $indent++;
        # 'foreach' provides a list context.
        foreach my $val ($obj->value()) {
            dump_object($val, $indent);
        }
        printf("%sENDGROUP\n", $istr);
    }
}
```

EXAMPLE 8 Creating a New Group Record and Writing to a File

```
# Prototype list of catalog tags and values.
my @items = (
    [ &EXT_STRING | &EXC_DEFAULT | &EXD_CREATOR      => "me"          ],
    [ &EXT_UINT32 | &EXC_DEFAULT | &EXD_PROC_PID     => $$              ],
    [ &EXT_UINT32 | &EXC_DEFAULT | &EXD_PROC_UID     => $<              ],
    [ &EXT_UINT32 | &EXC_DEFAULT | &EXD_PROC_GID     => $(              ],
    [ &EXT_STRING | &EXC_DEFAULT | &EXD_PROC_COMMAND => "/bin/stuff" ],
);

# Create a new group catalog object.
my $cat = new_catalog(&EXT_GROUP | &EXC_DEFAULT | &EXD_NONE);

# Create a new Group object and retrieve its data array.
my $group = new_group($cat);
my $ary = $group->value();

# Push the new Items onto the Group array.
foreach my $v (@items) {
    push(@$ary, new_item(new_catalog($v->[0]), $v->[1]));
}

# Nest the group within itself (performs a deep copy).
push(@$ary, $group);

# Dump out the group.
dump_object($group);
```

EXAMPLE 9 Dumping an exact File

```
#!/usr/bin/perl

use strict;
```

```

use warnings;
use blib;
use Sun::Solaris::Exacct qw(:EXACCT_ALL);

die("Usage is dumpexacct

# Open the exact file and display the header information.
my $ef = ea_new_file($ARGV[0], &O_RDONLY) || die(error_str());
printf("Creator: %s\n", $ef->creator());
printf("Hostname: %s\n\n", $ef->hostname());

# Dump the file contents
while (my $obj = $ef->get()) {
    ea_dump_object($obj);
}

# Report any errors
if (ea_error() != EXR_OK && ea_error() != EXR_EOF) {
    printf("\nERROR: %s\n", ea_error_str());
    exit(1);
}
exit(0);

```

Output From the Perl dump Method

This example shows the formatted output of the `Sun::Solaris::Exacct::Object->dump()` method.

```

GROUP
  Catalog = EXT_GROUP|EXC_DEFAULT|EXD_GROUP_PROC_PARTIAL
  ITEM
    Catalog = EXT_UINT32|EXC_DEFAULT|EXD_PROC_PID
    Value = 3
  ITEM
    Catalog = EXT_UINT32|EXC_DEFAULT|EXD_PROC_UID
    Value = 0
  ITEM
    Catalog = EXT_UINT32|EXC_DEFAULT|EXD_PROC_GID
    Value = 0
  ITEM
    Catalog = EXT_UINT32|EXC_DEFAULT|EXD_PROC_PROJID
    Value = 0
  ITEM
    Catalog = EXT_UINT32|EXC_DEFAULT|EXD_PROC_TASKID
    Value = 0
  ITEM

```

```
    Catalog = EXT_UINT64|EXC_DEFAULT|EXD_PROC_CPU_USER_SEC
    Value = 0
ITEM
    Catalog = EXT_UINT64|EXC_DEFAULT|EXD_PROC_CPU_USER_NSEC
    Value = 0
ITEM
    Catalog = EXT_UINT64|EXC_DEFAULT|EXD_PROC_CPU_SYS_SEC
    Value = 890
ITEM
    Catalog = EXT_UINT64|EXC_DEFAULT|EXD_PROC_CPU_SYS_NSEC
    Value = 760000000
ITEM
    Catalog = EXT_UINT64|EXC_DEFAULT|EXD_PROC_START_SEC
    Value = 1011869897
ITEM
    Catalog = EXT_UINT64|EXC_DEFAULT|EXD_PROC_START_NSEC
    Value = 380771911
ITEM
    Catalog = EXT_UINT64|EXC_DEFAULT|EXD_PROC_FINISH_SEC
    Value = 0
ITEM
    Catalog = EXT_UINT64|EXC_DEFAULT|EXD_PROC_FINISH_NSEC
    Value = 0
ITEM
    Catalog = EXT_STRING|EXC_DEFAULT|EXD_PROC_COMMAND
    Value = fsflush
ITEM
    Catalog = EXT_UINT32|EXC_DEFAULT|EXD_PROC_TTY_MAJOR
    Value = 4294967295
ITEM
    Catalog = EXT_UINT32|EXC_DEFAULT|EXD_PROC_TTY_MINOR
    Value = 4294967295
ITEM
    Catalog = EXT_STRING|EXC_DEFAULT|EXD_PROC_HOSTNAME
    Value = mower
ITEM
    Catalog = EXT_UINT64|EXC_DEFAULT|EXD_PROC_FAULTS_MAJOR
    Value = 0
ITEM
    Catalog = EXT_UINT64|EXC_DEFAULT|EXD_PROC_FAULTS_MINOR
    Value = 0
ITEM
    Catalog = EXT_UINT64|EXC_DEFAULT|EXD_PROC_MESSAGES_SND
    Value = 0
ITEM
    Catalog = EXT_UINT64|EXC_DEFAULT|EXD_PROC_MESSAGES_RCV
    Value = 0
ITEM
```



```
Catalog = EXT_UINT64|EXC_DEFAULT|EXD_PROC_BLOCKS_IN
Value = 19
ITEM
Catalog = EXT_UINT64|EXC_DEFAULT|EXD_PROC_BLOCKS_OUT
Value = 40833
ITEM
Catalog = EXT_UINT64|EXC_DEFAULT|EXD_PROC_CHARS_RDWR
Value = 0
ITEM
Catalog = EXT_UINT64|EXC_DEFAULT|EXD_PROC_CONTEXT_VOL
Value = 129747
ITEM
Catalog = EXT_UINT64|EXC_DEFAULT|EXD_PROC_CONTEXT_INV
Value = 79
ITEM
Catalog = EXT_UINT64|EXC_DEFAULT|EXD_PROC_SIGNALS
Value = 0
ITEM
Catalog = EXT_UINT64|EXC_DEFAULT|EXD_PROC_SYSCALLS
Value = 0
ITEM
Catalog = EXT_UINT32|EXC_DEFAULT|EXD_PROC_ACCT_FLAGS
Value = 1
ITEM
Catalog = EXT_UINT32|EXC_DEFAULT|EXD_PROC_ANCPID
Value = 0
ITEM
Catalog = EXT_UINT32|EXC_DEFAULT|EXD_PROC_WAIT_STATUS
Value = 0
ENDGROUP
```


Resource Controls in Oracle Solaris

This chapter describes resource controls and their properties.

- [“Overview of Resource Controls” on page 51](#)
- [“Resource Controls Flags and Actions” on page 52](#)
- [“Resource Controls API Functions” on page 63](#)
- [“Resource Control Code Examples” on page 64](#)
- [“Programming Issues Associated With Resource Controls” on page 69](#)

Overview of Resource Controls

Use the extended accounting facility to determine the resource consumption of workloads on your system. After the resource consumption has been determined, use the resource control facility to place bounds on resource usage. Bounds that are placed on resources prevent workloads from over-consuming resources.

For an overview of resource controls and example commands for administering resource controls, see [Chapter 6, “About Resource Controls” in *Administering Resource Management in Oracle Solaris 11.3*](#) and [Chapter 7, “Administering Resource Controls Tasks” in *Administering Resource Management in Oracle Solaris 11.3*](#).

The resource control facility adds the following benefits.

- **Dynamically set**
Resource controls can be adjusted while the system is running.
- **Containment level granularity**
Resource controls are arranged in a containment level of zone, project, task, or process. The containment level simplifies the configuration and aligns the collected values closer to the particular zone, project, task, or process.

Resource Controls Flags and Actions

This section describes flags, actions, and signals associated with resource controls.

rlimit, Resource Limit

`rlimit` is process-based. `rlimit` establishes a restricting boundary on the consumption of a variety of system resources by a process. Each process that the process creates inherits from the original process. A resource limit is defined by a pair of values. The values specify the current (soft) limit and the maximum (hard) limit.

A process might irreversibly lower its hard limit to any value that is greater than or equal to the soft limit. Only a process with root ID can raise the hard limit. See `setrlimit()` and `getrlimit()`.

The `rlimit` structure contains two members that define the soft limit and hard limit.

```
rlim_t    rlim_cur;    /* current (soft) limit */
rlim_t    rlim_max    /* hard limit */
```

rctl, Resource Control

`rctl` extends the process-based limits of `rlimit` by controlling resource consumption by processes, tasks, and projects defined in the project database.

Note - The `rctl` mechanism is preferred to the use of `rlimit` to set resource limits. The only reason to use the `rlimit` facility is when portability is required across UNIX platforms.

Applications fall into the following broad categories depending on how an application deals with resource controls. Based on the action that is taken, resource controls can be further classified. Most report an error and terminate operation. Other resource controls allow applications to resume operation and adapt to the reduced resource usage. A progressive chain of actions at increasing values can be specified for each resource control.

The list of attributes for a resource control consists of a privilege level, a threshold value, and an action that is taken when the threshold is exceeded.

Resource Control Values and Privilege Levels

Each threshold value on a resource control must be associated with one of the following privilege levels:

RCPRIV_BASIC

Privilege level can be modified by the owner of the calling process. RCPRIV_BASIC is associated with a resource's soft limit.

RCPRIV_PRIVILEGED

Privilege level can be modified only by privileged (root) callers. RCPRIV_PRIVILEGED is associated with a resource's hard limit.

`set rctl(2)` will only succeed when called as a privileged user in the global zone. Inside a non-global zone, root cannot set zone-wide controls.

RCPRIV_SYSTEM

Privilege level remains fixed for the duration of the operating system instance.

[Figure 4, “Setting Privilege Levels for Signals,” on page 63](#) shows the timeline for setting privilege levels for signals that are defined by the `/etc/project` file `process.max-cpu-time` resource control.

Local Actions and Local Flags

The local action and local flags are applied to the current resource control value represented by this resource control block. Local actions and local flags are value-specific. For each threshold value that is placed on a resource control, the following local actions and local flags are available:

RCTL_LOCAL_NOACTION

No local action is taken when this resource control value is exceeded.

RCTL_LOCAL_SIGNAL

The specified signal, set by `rctlblk_set_local_action()`, is sent to the process that placed this resource control value in the value sequence.

RCTL_LOCAL_DENY

When this resource control value is encountered, the request for the resource is denied. Set on all values if `RCTL_GLOBAL_DENY_ALWAYS` is set for this control. Cleared on all values if `RCTL_GLOBAL_DENY_NEVER` is set for this control.

RCTL_LOCAL_MAXIMAL

This resource control value represents a request for the maximum amount of resource for this control. If RCTL_GLOBAL_INFINITE is set for this resource control, RCTL_LOCAL_MAXIMAL indicates an unlimited resource control value that is never exceeded.

Global Actions and Global Flags

Global flags apply to all current resource control values represented by this resource control block. Global actions and global flags are set by `rctladm(1M)`. Global actions and global flags cannot be set with `setrctl()`. Global flags apply to all resource controls. For each threshold value that is placed on a resource control, the following global actions and global flags are available:

RCTL_GLOBAL_NOACTION

No global action is taken when a resource control value is exceeded on this control.

RCTL_GLOBAL_SYSLOG

A standard message is logged by the `syslog()` facility when any resource control value on a sequence associated with this control is exceeded.

RCTL_GLOBAL_SECONDS

Defines the unit string of the limit value as seconds.

RCTL_GLOBAL_COUNT

Defines the unit string of the limit value as count.

RCTL_GLOBAL_BYTES

Defines the unit string of the limit value as bytes.

RCTL_GLOBAL_SYSLOG_NEVER

Flag means that RCTL_GLOBAL_SYSLOG cannot be set for this resource control through `rctladm(1M)`.

RCTL_GLOBAL_NOBASIC

No values with the RCPRIV_BASIC privilege are permitted on this control.

RCTL_GLOBAL_LOWERABLE

Non-privileged callers are able to lower the value of privileged resource control values on this control.

RCTL_GLOBAL_DENY_ALWAYS

The action that is taken when a control value is exceeded on this control always includes denial of the resource.

RCTL_GLOBAL_DENY_NEVER

The action that is taken when a control value is exceeded on this control always excludes denial of the resource. The resource is always granted, although other actions can also be taken.

RCTL_GLOBAL_FILE_SIZE

The valid signals for local actions include the SIGXFSZ signal.

RCTL_GLOBAL_CPU_TIME

The valid signals for local actions include the SIGXCPU signal.

RCTL_GLOBAL_SIGNAL_NEVER

No local actions are permitted on this control. The resource is always granted.

RCTL_GLOBAL_INFINITE

This resource control supports the concept of an unlimited value. Generally, an unlimited value applies only to accumulation-oriented resources, such as CPU time.

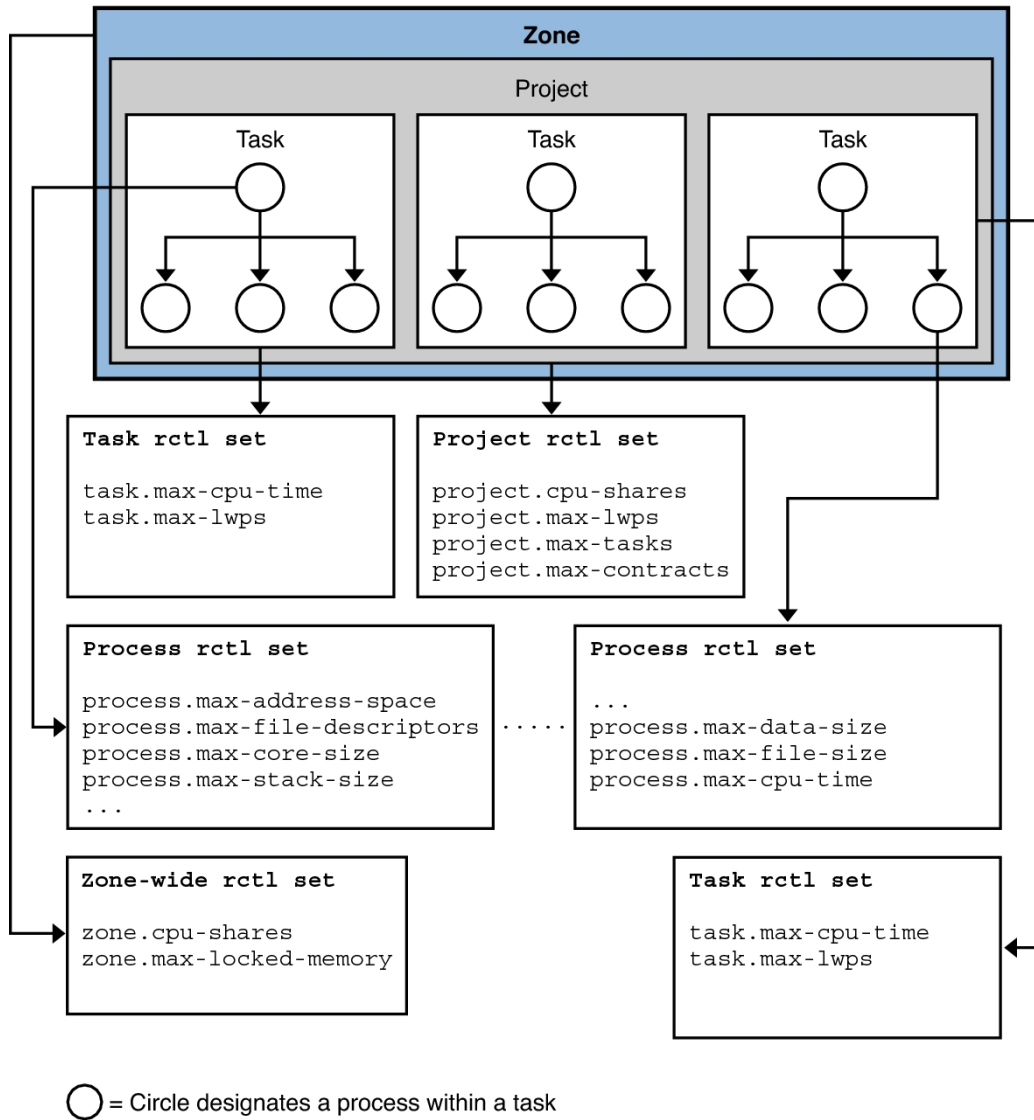
RCTL_GLOBAL_UNOBSERVABLE

Generally, a task or project related resource control does not support observational control values. An RCPRIV_BASIC privileged control value placed on a task or process generates an action only if the value is exceeded by the process that placed the value.

Resource Control Sets Associated With a Zone, Project, Processes, and Tasks

The following figure shows the resource control sets associated with zones, tasks, processes and a project.

FIGURE 3 Resource Control Sets for Zone, Task, Project, and Process



More than one resource control can exist on a resource, each resource control at a containment level in the process model. Resource controls can be active on the same resource for both a

process and collective task or collective project. In this case, the action for the process takes precedence. For example, action is taken on `process.max-cpu-time` before `task.max-cpu-time` if both controls are encountered simultaneously.

Resource Controls Associated With a Project

Resource controls associated with a project include the following:

`project.cpu-cap`

Absolute limit on the amount of CPU resources that can be consumed by a project. A value of 100 means 100 percent of one CPU as the `project.cpu-cap` setting. A value of 125 is 125 percent, because 100 percent corresponds to one full CPU on the system when using CPU caps.

`project.cpu-shares`

The number of CPU shares that are granted to this project for use with the fair share scheduler, FSS.

`project.max-crypto-memory`

Total amount of kernel memory that can be used by `libpkcs11` for hardware crypto acceleration. Allocations for kernel buffers and session-related structures are charged against this resource control.

`project.max-locked-memory`

Total amount of physical locked memory allowed.

Note that this resource control replaced `project.max-device-locked-memory`, which has been removed.

`project.max-msg-ids`

Maximum number of System V message queues allowed for a project.

`project.max-port-ids`

Maximum allowable number of event ports.

`project.max-processes`

Maximum number of process table slots simultaneously available to this project.

Note - Both normal processes and zombie processes take up process table slots. The `max-processes` resource control thus protects against zombie processes exhausting the process table. Note that `max-lwps` cannot protect against zombie processes exhausting the process table since zombie processes do not have any LWPs by definition.

`project.max-sem-ids`

Maximum number of semaphore IDs allowed for a project.

`project.max-shm-ids`

Maximum number of shared memory IDs allowed for this project.

`project.max-msg-ids`

Maximum number of message queue IDs allowed for this project.

`project.max-shm-memory`

Total amount of System V shared memory allowed for this project.

`project.max-lwps`

Maximum number of LWPs simultaneously available to this project.

`project.max-tasks`

Maximum number of tasks allowable in this project.

`project.max-contracts`

Maximum number of contracts allowed in this project.

Resource Controls Associated With Tasks

Resource controls associated with tasks include the following:

`task.max-cpu-time`

Maximum CPU time (seconds) available to this task's processes.

`task.max-lwps`

Maximum number of LWPs simultaneously available to this task's processes.

`task.max-processes`

Maximum number of process table slots simultaneously available to this task's processes.

Note - Both normal processes and zombie processes take up process table slots. The `max-processes` resource control thus protects against zombie processes exhausting the process table. Note that `max-lwps` cannot protect against zombie processes exhausting the process table since zombie processes do not have any LWPs by definition.

Resource Controls Associated With Processes

Resource controls associated with processes include the following:

`process.max-address-space`

Maximum amount of address space (bytes), as summed over segment sizes, available to this process.

`process.max-core-size`

Maximum size (bytes) of a core file that is created by this process.

`process.max-cpu-time`

Maximum CPU time (seconds) available to this process.

`process.max-file-descriptor`

Maximum file descriptor index that is available to this process.

`process.max-file-size`

Maximum file offset (bytes) available for writing by this process.

`process.max-msg-messages`

Maximum number of messages on a message queue. This value is copied from the resource control at `msgget()` time.

`process.max-msg-qbytes`

Maximum number (bytes) of messages on a message queue. This value is copied from the resource control at `msgget()` time. When you set a new `project.max-msg-qbytes` value, initialization occurs only on the subsequently created values. The new `project.max-msg-qbytes` value does not effect existing values.

`process.max-sem-nsems`

Maximum number of semaphores allowed for a semaphore set.

`process.max-sem-ops`

Maximum number of semaphore operations that are allowed for a `semop()` call. This value is copied from the resource control at `msgget()` time. A new `project.max-sem-ops` value only affects the initialization of subsequently created values and has no effect on existing values.

`process.max-port-events`

Maximum number of events that are allowed per event port.

Zone-Wide Resource Controls

Zone-wide resource controls are available on a system with zones installed. Zone-wide resource controls limit the total resource usage of all process entities within a zone.

`zone.cpu-cap`

Absolute limit on the amount of CPU resources that can be consumed by a non-global zone. A value of `100` means 100 percent of one CPU as the `project.cpu-cap` setting. A value of `125` is 125 percent, because 100 percent corresponds to one full CPU on the system when using CPU caps.

`zone.cpu-shares`

Limit on the number of fair share scheduler (FSS) CPU shares for a zone. The scheduling class must be FSS. CPU shares are first allocated to the zone, and then further subdivided among projects within the zone as specified in the `project.cpu-shares` entries. A zone with a higher number of `zone.cpu-shares` is allowed to use more CPU than a zone with a low number of shares.

`zone.max-locked-memory`

Total amount of physical locked memory available to a zone.

`zone.max-lofi`

Maximum number of `lofi` devices that can be created by a zone.

`zone.max-lwps`

Maximum number of LWPs simultaneously available to this zone

`zone.max-msg-ids`

Maximum number of message queue IDs allowed for this zone

`zone.max-processes`

Maximum number of process table slots simultaneously available to this zone

Note - The `zone.max-processes` resource control enhances resource isolation by preventing a zone from using too many process table slots and thus affecting other zones. The allocation of the process table slots resource across projects within the zone can be controlled by using the `project.max-processes` resource control. The global property name for this control is `max-processes`. The `zone.max-processes` resource control can also encompass the `zone.max-lwps` resource control. If `zone.max-processes` is set and `zone.max-lwps` is not set, then `zone.max-lwps` is implicitly set to 10 times the `zone.max-processes` when the zone is booted.

Both normal processes and zombie processes take up process table slots. The `max-processes` resource control thus protects against zombie processes exhausting the process table. Note that `max-lwps` cannot protect against zombie processes exhausting the process table since zombie processes do not have any LWPs by definition.

`zone.max-sem-ids`

Maximum number of semaphore IDs allowed for this zone

`zone.max-shm-ids`

Maximum number of shared memory IDs allowed for this zone

`zone.max-shm-memory`

Total amount of shared memory allowed for this zone

`zone.max-swap`

Total amount of swap that can be consumed by user process address space mappings and `tmpfs` mounts for this zone.

For more information, see [“Zone-Wide Resource Controls”](#) in *Administering Resource Management in Oracle Solaris 11.3*.

For information on configuring zone-wide resource controls, see [Chapter 1, “Configuration Resources for Non-Global Zones”](#) in *Oracle Solaris Zones Configuration Resources* and [Chapter 1, “How to Plan and Configure Non-Global Zones”](#) in *Creating and Using Oracle Solaris Zones*.

Note that it is possible to use the `zonecfg` command to apply a zone-wide resource control to the global zone on a system with non-global zones installed. Also note that the `setctrl` command will only succeed when called as a privileged user in the global zone. Inside a non-global zone, `root` cannot set zone-wide resource controls.

Signals Used With Resource Controls

For each threshold value that is placed on a resource control, the following restricted set of signals is available:

SIGBART

Terminate the process.

SIGXRES

Signal generated by resource control facility when the resource control limit is exceeded.

SIGHUP

When carrier drops on an open line, the process group that controls the terminal is sent a hangup signal, SIGHUP.

SIGSTOP

Job control signal. Stop the process. Stop signal not from terminal.

SIGTERM

Terminate the process. Termination signal sent by software.

SIGKILL

Terminate the process. Kill the program.

SIGXFSX

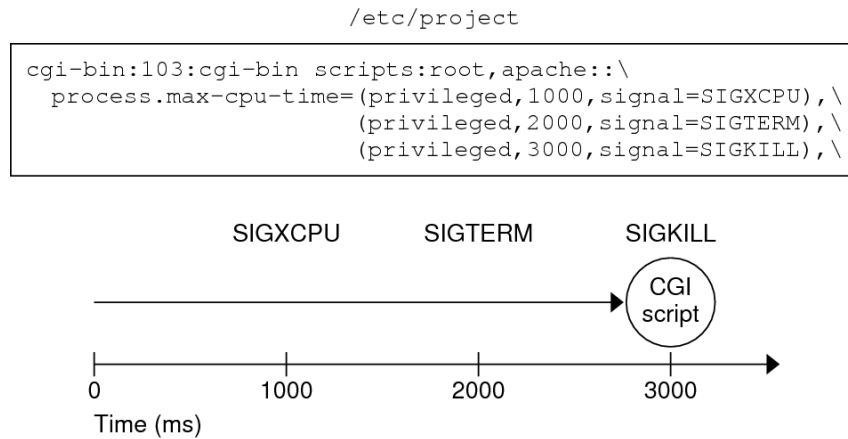
Terminate the process. File size limit exceeded. Available only to resource controls with the RCTL_GLOBAL_FILE_SIZE property.

SIGXCPU

Terminate the process. CPU time limit exceeded. Available only to resource controls with the RCTL_GLOBAL_CPU_TIME property.

Other signals might be permitted due to global properties of a specific control.

Note - Calls to `setrctl()` with illegal signals fail.

FIGURE 4 Setting Privilege Levels for Signals

Resource Controls API Functions

The resource controls API contains functions that:

- [“Operate on Action-Value Pairs of a Resource Control” on page 63](#)
- [“Operate on Local Modifiable Values” on page 64](#)
- [“Retrieve Local Read-Only Values” on page 64](#)
- [“Retrieve Global Read-Only Actions” on page 64](#)

Operate on Action-Value Pairs of a Resource Control

The following list contains the functions that set or get the resource control block.

[setrctl\(2\)](#)

[getrctl\(2\)](#)

Operate on Local Modifiable Values

The following list contains the functions associated with the local, modifiable resource control block.

<code>rctlblk_set_privilege(3C)</code>	<code>rctlblk_set_local_action(3C)</code>
<code>rctlblk_get_privilege(3C)</code>	<code>rctlblk_get_local_action(3C)</code>
<code>rctlblk_set_value(3C)</code>	<code>rctlblk_set_local_flags(3C)</code>
<code>rctlblk_get_value(3C)</code>	<code>rctlblk_get_local_flags(3C)</code>

Retrieve Local Read-Only Values

The following list contains the functions associated with the local, read-only resource control block.

```
rctlblk_get_recipient_pid(3C)
rctlblk_get_firing_time(3C)
rctlblk_get_enforced_value(3C)
```

Retrieve Global Read-Only Actions

The following list contains the functions associated with the global, read-only resource control block.

```
rctlblk_get_global_action(3C)
rctlblk_get_global_flags(3C)
```

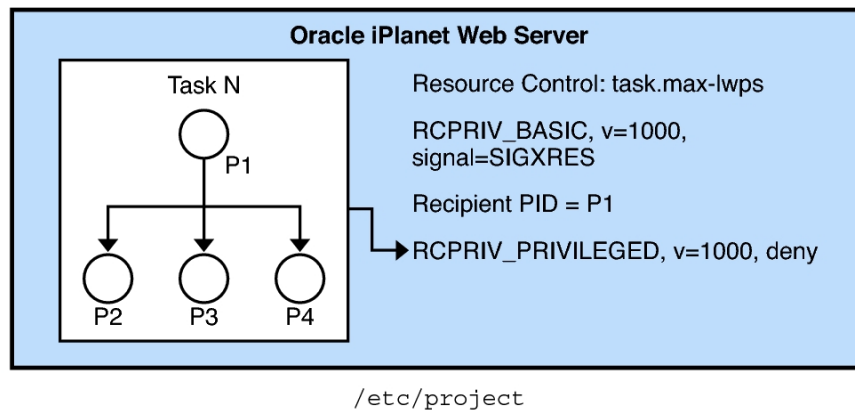
Resource Control Code Examples

Master Observing Process for Resource Controls

The following example is the master observer process. [Figure 5, “Master Observing Process,” on page 65](#) shows the resource controls for the master observing process.

Note - The line break is not valid in an `/etc/project` file. The line break is shown here only to allow the example to display on a printed or displayed page. Each entry in the `/etc/project` file must be on a separate line.

FIGURE 5 Master Observing Process



```
...
OracleiPlanetWebServer:200:Oracle iPlanet Web Application Server:root::\
task.max=lwps=(PRIVILEGED,1000,deny)
```

The key points for the example include the following:

- Because the task's limit is privileged, the application cannot change the limit, or specify an action, such as a signal. A master process solves this problem by establishing the same resource control as a basic resource control on the task. The master process uses the same value or a little less on the resource, but with a different action, `signal = XRES`. The master process creates a thread to wait for this signal.
- The `rctlblk` is opaque. The struct needs to be dynamically allocated.
- Note the blocking of all signals before creating the thread, as required by `sigwait(2)`.
- The thread calls `sigwait(2)` to block for the signal. If `sigwait()` returns the `SIGXRES` signal, the thread notifies the master process' children, which adapts to reduce the number of LWPs being used. Each child should also be modelled similarly, with a thread in each child, waiting for this signal, and adapting its process' LWP usage appropriately.

```
rctlblk_t *mlwprcb;
sigset_t smask;
```

```
/* Omit return value checking/error processing to keep code sample short */
/* First, install a RCPRIV_BASIC, v=1000, signal=SIGXRES rctl */
mlwprcb = calloc(1, rctlblk_size()); /* rctl blocks are opaque: */
    rctlblk_set_value(mlwprcb, 1000);
    rctlblk_set_privilege(mlwprcb, RCPRIV_BASIC);
    rctlblk_set_local_action(mlwprcb, RCTL_LOCAL_SIGNAL, SIGXRES);
    if (setrctl("task.max-lwps", NULL, mlwprcb, RCTL_INSERT) == -1) {
        perror("setrctl");
        exit (1);
    }

/* Now, create the thread which waits for the signal */
    sigemptyset(&smask);
    sigaddset(&smask, SIGXRES);
    thr_sigsetmask(SIG_BLOCK, &smask, NULL);
thr_create(NULL, 0, sigthread, (void *)SIGXRES, THR_DETACHED, NULL));

/* Omit return value checking/error processing to keep code sample short */

void *sigthread(void *a)
{
    int sig = (int)a;
    int rsig;
    sigset_t sset;

    sigemptyset(&sset);
    sigaddset(&sset, sig);

    while (1) {
        rsig = sigwait(&sset);
        if (rsig == SIGXRES) {
            notify_all_children();
            /* e.g. sigsend(P_PID, child_pid, SIGXRES); */
        }
    }
}
```

List all the Value-Action Pairs for a Specific Resource Control

The following example lists all the value-action pairs for a specific resource control, `task.max-lwps`. The key point for the example is that `getrctl(2)` takes two resource control blocks, and returns the resource control block for the `RCTL_NEXT` flag. To iterate through all resource

control blocks, repeatedly swap the resource control block values, as shown here using the `rcb_tmp` `rctl` block.

```

rctlblk_t *rcb1, *rcb2, *rcb_tmp;
...
/* Omit return value checking/error processing to keep code sample short */
rcb1 = calloc(1, rctlblk_size()); /* rctl blocks are opaque: */
                                /* "rctlblk_t rcb" does not work */
rcb2 = calloc(1, rctlblk_size());
getrctl("task.max-lwps", NULL, rcb1, RCTL_FIRST);
while (1) {
    print_rctl(rcb1);
    rcb_tmp = rcb2;
    rcb2 = rcb1;
    rcb1 = rcb_tmp;          /* swap rcb1 with rcb2 */
    if (getrctl("task.max-lwps", rcb2, rcb1, RCTL_NEXT) == -1) {
        if (errno == ENOENT) {
            break;
        }
    } else {
        perror("getrctl");
        exit (1);
    }
}
}

```

Set `project.cpu-shares` and Add a New Value

The key points of the example include the following:

- This example is similar to the example shown in [“Set `pool.comment` Property and Add New Property” on page 85](#).
- Use `bcopy()`, rather than buffer swapping as in [“List all the Value-Action Pairs for a Specific Resource Control” on page 66](#).
- To change the resource control value, call `setrctl()` with the `RCTL_REPLACE` flag. The new resource control block is identical to the old resource control block except for the new control value.

```

rctlblk_set_value(blk1, nshares);
if (setrctl("project.cpu-shares", blk2, blk1, RCTL_REPLACE) != 0)

```

The example gets the project's CPU share allocation, `project.cpu-shares`, and changes its value to `nshares`.

```

/* Omit return value checking/error processing to keep code sample short */

```

```
blk1 = malloc(rctlblk_size());
getrctl("project.cpu-shares", NULL, blk1, RCTL_FIRST);
my_shares = rctlblk_get_value(blk1);
printout_my_shares(my_shares);
/* if privileged, do the following to */
/* change project.cpu-shares to "nshares" */
blk1 = malloc(rctlblk_size());
blk2 = malloc(rctlblk_size());
if (getrctl("project.cpu-shares", NULL, blk1, RCTL_FIRST) != 0) {
    perror("getrctl failed");
    exit(1);
}
bcopy(blk1, blk2, rctlblk_size());
rctlblk_set_value(blk1, nshares);
if (setrctl("project.cpu-shares", blk2, blk1, RCTL_REPLACE) != 0) {
    perror("setrctl failed");
    exit(1);
}
```

Set LWP Limit Using Resource Control Blocks

In the following example, an application has set a privileged limit of 3000 LWPs that may not be exceeded. In addition, the application has set a basic limit of 2000 LWPs. When this limit is exceeded, a SIGXRES is sent to the application. Upon receiving a SIGXRES, the application might send notification to its child processes that might in turn reduce the number of LWPs the processes use or need.

```
/* Omit return value and error checking */

#include <rctl.h>

rctlblk_t *rcb1, *rcb2;

/*
 * Resource control blocks are opaque
 * and must be explicitly allocated.
 */
rcb1 = calloc(rctlblk_size());

rcb2 = calloc(rctlblk_size());

/* Install an RCPRIV_PRIVILEGED, v=3000: do not allow more than 3000 LWPs */
rctlblk_set_value(rcb1, 3000);
rctlblk_set_privilege(rcb1, RCPRIV_PRIVILEGED);
rctlblk_set_local_action(rcb1, RCTL_LOCAL_DENY);
```

```
setrctl("task.max-lwps", NULL, rcb1, RCTL_INSERT);

/* Install an RCPRIV_BASIC, v=2000 to send SIGXRES when LwPs exceeds 2000 */
rctlblk_set_value(rcb2, 2000);
rctlblk_set_privilege(rcb2, RCPRIV_BASIC);
rctlblk_set_local_action(rcb2, RCTL_LOCAL_SIGNAL, SIGXRES);
setrctl("task.max-lwps", NULL, rcb2, RCTL_INSERT);
```

Programming Issues Associated With Resource Controls

Consider the following issues when writing your application:

- The resource control block is opaque. The control block needs to be dynamically allocated.
- If a basic resource control is established on a task or project, the process that establishes this resource control becomes an observer. The action for this resource control block is applied to the observer. However, some resources cannot be observed in this manner.
- If a privileged resource control is set on a task or project, no observer process exists. However, any process that violates the limit becomes the subject of the resource control action.
- Only one action is permitted for each type: global and local.
- Only one basic `rctl` is allowed per process per resource control.

zonestat Utility for Monitoring Zones Resource Usage

The `zonestat` utility reports on the CPU, memory, and resource control utilization of the currently running zones. Each zone's utilization is reported as a percentage of both system resources and the zone's configured limits. For more information, see [Chapter 7, “Design Considerations for Resource Management Applications in Oracle Solaris Zones”](#), the `zonestat(1)` man page, and [“Reporting Active Zone Statistics with the zonestat Utility” in *Creating and Using Oracle Solaris Zones*](#).

Resource Pools in Oracle Solaris

This chapter describes resource pools and their properties.

- [“Overview of Resource Pools” on page 71](#)
- [“Dynamic Resource Pool Constraints and Objectives” on page 72](#)
- [“Resource Pools API Functions” on page 77](#)
- [“Resource Pool Code Examples” on page 83](#)
- [“Programming Issues Associated With Resource Pools” on page 86](#)

Overview of Resource Pools

Resource pools provide a framework for managing processor sets and thread scheduling classes. Resource pools are used for partitioning system resources. Resource pools enable you to separate workloads so that workload consumption of certain resources does not overlap. The resource reservation helps to achieve predictable performance on systems with mixed workloads.

For an overview of resource pools and example commands for administering resource pools, see [Chapter 12, “About Resource Pools” in *Administering Resource Management in Oracle Solaris 11.3*](#) and [Chapter 13, “Creating and Administering Resource Pools Tasks” in *Administering Resource Management in Oracle Solaris 11.3*](#).

A processor set groups the CPUs on a system into a bounded entity, on which a process or processes can run exclusively. Processes cannot extend beyond the processor set, nor can other processes extend into the processor set. A processor set enables tasks of similar characteristics to be grouped together and a hard upper boundary for CPU use to be set.

The resource pool framework allows the definition of a soft processor set with a maximum and minimum CPU count requirement. Additionally, the framework provides a hard-defined scheduling class for that processor set.

A zone can be bound to a resource pool through the `pool` property of the zone configuration. The zone is bound to the specified pool upon creation of the zone. The pool configuration can be changed only from the global zone. Zones cannot span multiple pools. All processes in a zone run in the same pool. However, multiple zones can bind to the same resource pool.

A resource pool defines:

- Processor set groups
- Scheduling class

Scheduling Class

Scheduling classes provide different CPU access characteristics to threads that are based on algorithmic logic. The scheduling classes include:

- Realtime scheduling class
- Interactive scheduling class
- Fixed priority scheduling class
- Timesharing scheduling class
- Fair share scheduling class

For an overview of fair share scheduler and example commands for administering the fair share scheduler, see [Chapter 8, “About Fair Share Scheduler” in *Administering Resource Management in Oracle Solaris 11.3*](#) and [Chapter 9, “Administering the Fair Share Scheduler Tasks” in *Administering Resource Management in Oracle Solaris 11.3*](#).

Do not mix scheduling classes in a set of CPUs. If scheduling classes are mixed in a CPU set, system performance might become erratic and unpredictable. Use processor sets to segregate applications by their characteristics. Assign scheduling classes under which the application best performs. For more information about the characteristics of an individual scheduling class, see `pricntl(1)`.

For an overview of resource pools and a discussion of when to use pools, see [Chapter 6, “Resource Pools in Oracle Solaris”](#).

Dynamic Resource Pool Constraints and Objectives

The `libpool` library defines properties that are available to the various entities that are managed within the pools facility. Each property falls into the following categories:

Configuration constraints

A constraint defines boundaries of a property. Typical constraints are the maximum and minimum allocations specified in the `libpool` configuration.

Objective

An objective changes the resource assignments of the current configuration to generate new candidate configurations that observe the established constraints. An objective has the following categories:

Workload dependent

A workload-dependent objective varies according to the conditions imposed by the workload. An example of the workload dependent objective is the `utilization` objective.

Workload independent

A workload-independent objective does not vary according to the conditions imposed by the workload. An example of the workload independent objective is the `CPU locality` objective.

An objective can take an optional prefix to indicate the importance of the objective. The objective is multiplied by this prefix, which is an integer from 0 to `INT64_MAX`, to determine the significance of the objective.

For usage examples, see [“How to Set Configuration Constraints” in *Administering Resource Management in Oracle Solaris 11.3*](#) and [“How to Define Configuration Objectives” in *Administering Resource Management in Oracle Solaris 11.3*](#).

System Properties

`system.bind-default` (writable boolean)

If the specified pool is not found in `/etc/project`, bind to pool with the `pool.default` property set to `TRUE`.

`system.comment` (writable string)

User description of system. `system.comment` is not used by the default pools commands, except when a configuration is initiated by the `poolcfg` utility. In this case, the system puts an informative message in the `system.comment` property for that configuration.

`system.name` (writable string)

User name for the configuration.

`system.version` (read-only integer)

libpool version required to manipulate this configuration.

Pools Properties

All pools properties except `pool.default` and `pool.sys_id` are writable.

`pool.active` (writable boolean)

If TRUE, mark this pool as active.

`pool.comment` (writable string)

User description of pool.

`pool.default` (read-only boolean)

If TRUE, mark this pool as the default pool. See the `system.bind-default` property.

`pool.importance` (writable integer)

Relative importance of this pool. Used for possible resource dispute resolution.

`pool.name` (writable string)

User name for pool. `setproject(3PROJECT)` uses `pool.name` as the value for the `project.pool` project attribute in the `project(4)` database.

`pool.scheduler` (writable string)

Scheduler class to which consumers of this pool are bound. This property is optional and if not specified, the scheduler bindings for consumers of this pool are not affected. For more information about the characteristics of an individual scheduling class, see `prioctl(1)`.

Scheduler classes include:

- RT for realtime scheduler
- TS for timesharing scheduler
- IA for interactive scheduler
- FSS for fair share scheduler
- FX for fixed priority scheduler

`pool.sys_id` (read-only integer)

This is the system-assigned pool ID.

Processor Set Properties

`pset.comment` (writable string)

User description of resource.

`pset.default` (read-only boolean)

Identifies the default processor set.

`pset.load` (read-only unsigned integer)

The load for this processor set. The lowest value is 0. The value increases in a linear fashion with the load on the set, as measured by the number of jobs in the system run queue.

`pset.max` (writable unsigned integer)

Maximum number of CPUs that are permitted in this processor set.

`pset.min` (writable unsigned integer)

Minimum number of CPUs that are permitted in this processor set.

`pset.name` (writable string)

User name for the resource.

`pset.size` (read-only unsigned integer)

Current number of CPUs in this processor set.

`pset.sys_id` (read-only integer)

System-assigned processor set ID.

`pset.type` (read-only string)

Names the resource type. Value for all processor sets is `pset`.

`pset.units` (read-only string)

Identifies the meaning of size-related properties. The value for all processor sets is `population`.

`cpu.comment` (writable string)

User description of CPU.

`pset.resmin`

If `pset.policy == minmax` (writable unsigned integer)

The minimum number of CPUs for the processor set.

If `pset.policy == assigned` (read-only unsigned integer)

The number of assigned CPUs, cores, or sockets.

`pset.resmax`

If `pset.policy == minmax` (writable unsigned integer)

The maximum number of CPUs for the processor set.

If `pset.policy == assigned` (read-only unsigned integer)

The number of assigned CPUs, cores, or sockets.

`pset.resize` (read-only string)

The current number of CPUs, cores, or sockets in the processor set. For `pset.policy == minmax`, this is always the number of CPUs.

`pset.restype`

The type of CPU resources assigned: CPUs, cores, or sockets. For `pset.policy == minmax`, this is always CPUs.

`pset.reslist`

The IDs of the assigned CPUs, cores, or sockets. For `pset.policy == minmax`, this is always blank.

`pset.policy`

`minmax` or `assigned`, depending on whether the processor set is using quantity-based allocation or specific assignment.

Using `libpool` to Manipulate Pool Configurations

The `libpool(3LIB)` pool configuration library defines the interface for reading and writing pools configuration files. The library also defines the interface for committing an existing configuration to becoming the running operating system configuration. The `<pool.h>` header provides type and function declarations for all library services.

The resource pools facility brings together process-bindable resources into a common abstraction that is called a pool. Processor sets and other entities can be configured, grouped, and labelled in a persistent fashion. Workload components can be associated with a subset of a system's total resources. The `libpool(3LIB)` library provides a C language API for accessing the resource pools facility. The `pooladm(1M)`, `poolbind(1M)`, and `poolcfg(1M)` make the resource pools facility available through command invocations from a shell.

Manipulate psets

The following list contains the man page and description of functions associated with creating or destroying psets and manipulating psets.

[processor_bind\(2\)](#)

Bind an LWP (lightweight process) or set of LWPs to a specified processor.

[pset_assign\(2\)](#)

Assign a processor to a processor set.

[pset_bind\(2\)](#)

Bind one or more LWPs (lightweight processes) to a processor set.

[pset_create\(2\)](#)

Create an empty processor set that contains no processors.

[pset_destroy\(2\)](#)

Destroy a processor set and release the associated constituent processors and processes.

[pset_setattr\(2\)](#), [pset_getattr\(2\)](#)

Set or get processor set attributes.

Resource Pools API Functions

This section lists all of the resource pool functions. Each function has a link to the man page and a short description of the function's purpose. The functions are divided into two groups, depending on whether the function performs an action or a query:

- [“Functions for Operating on Resource Pools and Associated Elements” on page 78](#)
- [“Functions for Querying Resource Pools and Associated Elements” on page 80](#)

The imported interfaces for libpool for swap sets is identical to the ones defined in this document.

Functions for Operating on Resource Pools and Associated Elements

The following list contains the man pages and descriptions of interfaces for performing actions related to pools and the associated elements.

`pool_associate(3POOL)`

Associate a resource with a specified pool.

`pool_component_to_elem(3POOL)`

Convert specified component to the pool element type.

`pool_conf_alloc(3POOL)`

Create a pool configuration.

`pool_conf_close(3POOL)`

Close the specified pool configuration and release the associated resources.

`pool_conf_commit(3POOL)`

Commit changes made to the specified pool configuration to permanent storage.

`pool_conf_export(3POOL)`

Save the given configuration to the specified location.

`pool_conf_free(3POOL)`

Release a pool configuration.

`pool_conf_open(3POOL)`

Create a pool configuration at the specified location.

`pool_conf_remove(3POOL)`

Removes the permanent storage for the configuration.

`pool_conf_rollback(3POOL)`

Restore the configuration state to the state that is held in the pool configuration's permanent storage.

`pool_conf_to_elem(3POOL)`

Convert specified pool configuration to the pool element type.

`pool_conf_update(3POOL)`

Update the library snapshot of kernel state.

`pool_create(3POOL)`

Create a new pool with the default properties and with default resources for each type.

`pool_destroy(3POOL)`

Destroy the specified pool. The associated resources are not modified.

`pool_dissociate(3POOL)`

Remove the association between the given resource and pool.

`pool_put_property(3POOL)`

Set the named property on the element to the specified value.

`pool_resource_create(3POOL)`

Create a new resource with the specified name and type for the provided configuration.

`pool_resource_destroy(3POOL)`

Remove the specified resource from the configuration file.

`pool_resource_to_elem(3POOL)`

Convert specified pool resource to the pool element type.

`pool_resource_transfer(3POOL)`

Transfer basic units from the source resource to the target resource.

`pool_resource_xtransfer(3POOL)`

Transfer the specified components from the source resource to the target resource.

`pool_rm_property(3POOL)`

Remove the named property from the element.

`pool_set_binding(3POOL)`

Bind the specified processes to the resources that are associated with pool on the running system.

`pool_set_status(3POOL)`

Modify the current state of the pools facility.

`pool_to_elem(3POOL)`

Convert specified pool to the pool element type.

`pool_value_alloc(3POOL)`

Allocate and return an opaque container for a pool property value.

`pool_value_free(3POOL)`

Release an allocated property values.

`pool_value_set_bool(3POOL)`

Set a property value of type boolean.

`pool_value_set_double(3POOL)`

Set a property value of type double.

`pool_value_set_int64(3POOL)`

Set a property value of type int64.

`pool_value_set_name(3POOL)`

Set a name=value pair for a pool property.

`pool_value_set_string(3POOL)`

Copy the string that was passed in.

`pool_value_set_uint64(3POOL)`

Set a property value of type uint64.

Functions for Querying Resource Pools and Associated Elements

The following list contains the man pages and descriptions of interfaces for performing queries related to pools and the associated elements.

`pool_component_info(3POOL)`

Return a string that describes the given component.

`pool_conf_info(3POOL)`

Return a string describing the entire configuration.

`pool_conf_location(3POOL)`

Return the location string that was provided to `pool_conf_open()` for the given specified configuration.

`pool_conf_status(3POOL)`

Return the validity status for a pool configuration.

`pool_conf_validate(3POOL)`

Check the validity of the contents of the given configuration.

`pool_dynamic_location(3POOL)`

Return the location that was used by the pools framework to store the dynamic configuration.

`pool_error(3POOL)`

Return the error value of the last failure that was recorded by calling a resource pool configuration library function.

`pool_get_binding(3POOL)`

Return the name of the pool on the running system that contains the set of resources to which the specified process is bound.

`pool_get_owing_resource(3POOL)`

Return the resource that currently contains the specified component.

`pool_get_pool(3POOL)`

Return the pool with the specified name from the provided configuration.

`pool_get_property(3POOL)`

Retrieve the value of the named property from the element.

`pool_get_resource(3POOL)`

Return the resource with the given name and type from the provided configuration.

`pool_get_resource_binding(3POOL)`

Return the name of the pool on the running system that contains the set of resources to which the given process is bound.

`pool_get_status(3POOL)`

Retrieve the current state of the pools facility.

`pool_info(3POOL)`

Return a description of the specified pool.

`pool_query_components(3POOL)`

Retrieve all resource components that match the specified list of properties.

`pool_query_pool_resources(3POOL)`

Return a null-terminated array of resources currently associated with the pool.

`pool_query_pools(3POOL)`

Return the list of pools that match the specified list of properties.

`pool_query_resource_components(3POOL)`

Return a null-terminated array of the components that make up the specified resource.

`pool_query_resources(3POOL)`

Return the list of resources that match the specified list of properties.

`pool_resource_info(3POOL)`

Return a description of the specified resource.

`pool_resource_type_list(3POOL)`

Enumerate the resource types that are supported by the pools framework on this platform.

`pool_static_location(3POOL)`

Return the location that was used by the pools framework to store the default configuration for pools framework instantiation.

`pool_strerror(3POOL)`

Return a description of each valid pool error code.

`pool_value_get_bool(3POOL)`

Get a property value of type `boolean`.

`pool_value_get_double(3POOL)`

Get a property value of type `double`.

`pool_value_get_int64(3POOL)`

Get a property value of type `int64`.

`pool_value_get_name(3POOL)`

Return the name that was assigned to the specified pool property.

`pool_value_get_string(3POOL)`

Get a property value of type `string`.

`pool_value_get_type(3POOL)`

Return the type of the data that is contained by the specified pool value.

`pool_value_get_uint64(3POOL)`

Get a property value of type `uint64`.

`pool_version(3POOL)`

Get the version number of the pool library.

`pool_walk_components(3POOL)`

Invoke callback on all components that are contained in the resource.

`pool_walk_pools(3POOL)`

Invoke callback on all pools that are defined in the configuration.

`pool_walk_properties(3POOL)`

Invoke callback on all properties defined for the given element.

`pool_walk_resources(3POOL)`

Invoke callback on all resources that are associated with the pool.

Resource Pool Code Examples

This section contains code examples of the resource pools interface.

Ascertain the Number of CPUs in the Resource Pool

`sysconf(3C)` provides information about the number of CPUs on an entire system. The following example provides the granularity of ascertaining the number of CPUs that are defined in a particular application's pools `pset`.

The key points for this example include the following:

- `pvals[]` should be a NULL terminated array.
- `pool_query_pool_resources()` returns a list of all resources that match the `pvals` array type `pset` from the application's pool `my_pool`. Because a pool can have only one instance of the `pset` resource, each instance is always returned in `nelem.reslist[]` contains only one element, the `pset` resource.

```
pool_value_t *pvals[2] = {NULL}; /* pvals[] should be NULL terminated */

/* NOTE: Return value checking/error processing omitted */
/* in all examples for brevity */

conf_loc = pool_dynamic_location();
conf = pool_conf_alloc();
pool_conf_open(conf, conf_loc, PO_RDONLY);
my_pool_name = pool_get_binding(getpid());
my_pool = pool_get_pool(conf, my_pool_name);
pvals[0] = pool_value_alloc();
pvals2[2] = { NULL, NULL };
pool_value_set_name(pvals[0], "type");
pool_value_set_string(pvals[0], "pset");

reslist = pool_query_pool_resources(conf, my_pool, &nelem, pvals);
pool_value_free(pvals[0]);
pool_query_resource_components(conf, reslist[0], &nelem, NULL);
printf("pool %s: %u cpu", my_pool_name, nelem);
pool_conf_close(conf);
```

List All Resource Pools

The following example lists all resource pools defined in an application's pools pset.

The key points of the example include the following:

- Open the dynamic conf file read-only, PO_RDONLY. `pool_query_pools()` returns the list of pools in `pl` and the number of pools in `nelem`. For each pool, call `pool_get_property()` to get the `pool.name` property from the element into the `pval` value.
- `pool_get_property()` calls `pool_to_elem()` to convert the `libpool` entity to an opaque value. `pool_value_get_string()` gets the string from the opaque pool value.

```
conf = pool_conf_alloc();
pool_conf_open(conf, pool_dynamic_location(), PO_RDONLY);
pl = pool_query_pools(conf, &nelem, NULL);
pval = pool_value_alloc();
for (i = 0; i < nelem; i++) {
    pool_get_property(conf, pool_to_elem(conf, pl[i]), "pool.name", pval);
    pool_value_get_string(pval, &fname);
    printf("%s\n", name);
}
pool_value_free(pval);
free(pl);
pool_conf_close(conf);
```

Report Pool Statistics for a Given Pool

The following example reports statistics for the designated pool.

The key points for the example include the following:

- `pool_query_pool_resources()` gets a list of all resources in `rl`. Because the last argument to `pool_query_pool_resources()` is `NULL`, all resources are returned. For each resource, the name, load and size properties are read, and printed.
- The call to `strdup()` allocates local memory and copies the string returned by `get_string()`. The call to `get_string()` returns a pointer that is freed by the next call to `get_property()`. If the call to `strdup()` is not included, subsequent references to the string(s) could cause the application to fail with a segmentation fault.

```
printf("pool %s\n:" pool_name);
pool = pool_get_pool(conf, pool_name);
rl = pool_query_pool_resources(conf, pool, &nelem, NULL);
for (i = 0; i < nelem; i++) {
    pool_get_property(conf, pool_resource_to_elem(conf, rl[i]), "type", pval);
    pool_value_get_string(pval, &type);
    type = strdup(type);
    snprintf(prop_name, 32, "%s.%s", type, "name");
    pool_get_property(conf, pool_resource_to_elem(conf, rl[i]),
        prop_name, pval);
    pool_value_get_string(val, &res_name);
    res_name = strdup(res_name);
    snprintf(prop_name, 32, "%s.%s", type, "load");
    pool_get_property(conf, pool_resource_to_elem(conf, rl[i]),
        prop_name, pval);
    pool_value_get_uint64(val, &load);
    snprintf(prop_name, 32, "%s.%s", type, "size");
    pool_get_property(conf, pool_resource_to_elem(conf, rl[i]),
        prop_name, pval);
    pool_value_get_uint64(val, &size);
    printf("resource %s: size %llu load %llu\n", res_name, size, load);
    free(type);
    free(res_name);
}
free(rl);
```

Set `pool.comment` Property and Add New Property

The following example sets the `pool.comment` property for the `pset`. The example also creates a new property in `pool.newprop`.

The key point for the example includes the following:

- In the call to `pool_conf_open()`, using `PO_RDWR` on a static configuration file, requires the caller to be root.
- To commit these changes to the `pset` after running this utility, issue a `pooladm -c` command. To have the utility commit the changes, call `pool_conf_commit()` with a nonzero second argument.

```
pool_set_comment(const char *pool_name, const char *comment)
{
    pool_t *pool;
    pool_elem_t *pool_elem;
    pool_value_t *pval = pool_value_alloc();
    pool_conf_t *conf = pool_conf_alloc();
    /* NOTE: need to be root to use PO_RDWR on static configuration file */
    pool_conf_open(conf, pool_static_location(), PO_RDWR);
    pool = pool_get_pool(conf, pool_name);
    pool_value_set_string(pval, comment);
    pool_elem = pool_to_elem(conf, pool);
    pool_put_property(conf, pool_elem, "pool.comment", pval);
    printf("pool %s: pool.comment set to %s\n:" pool_name, comment);
    /* Now, create a new property, customized to installation site */
    pool_value_set_string(pval, "New String Property");
    pool_put_property(conf, pool_elem, "pool.newprop", pval);
    pool_conf_commit(conf, 0); /* NOTE: use 0 to ensure only */
                               /* static file gets updated */

    pool_value_free(pval);
    pool_conf_close(conf);
    pool_conf_free(conf);
    /* NOTE: Use "pooladm -c" later, or pool_conf_commit(conf, 1) */
    /* above for changes to the running system */
}
```

An alternative way of modifying a pool's comment and adding a new pool property is to use `poolcfg(1M)`.

```
poolcfg -c 'modify pool pool-name (string pool.comment = "cmt-string")'
poolcfg -c 'modify pool pool-name (string pool.newprop =
            "New String Property")'
```

Programming Issues Associated With Resource Pools

Consider the following issues when writing your application.

- Each site can add its own list of properties to the pools configuration.

Multiple configurations can be maintained in multiple configuration files. The system administrator can commit different files to reflect changes to the resource consumption at different time slots. These time slots can include different times of the day, week, month, or seasons depending on load conditions.

- Resource sets can be shared between pools, but a pool has only one resource set of a given type. So, the `pset_default` can be shared between the default and a particular application's database pools.
- Use `pool_value_*`() interfaces carefully. Keep in mind the memory allocation issues for string pool values. See [“Report Pool Statistics for a Given Pool” on page 85](#).

zonestat Utility for Monitoring Resource Pools in Oracle Solaris Zones

The `zonestat` utility can be used to report on the CPU, memory, and resource control utilization of the currently running zones. Each zone's utilization is reported as a percentage of both system resources and the zone's configured limits. For more information, see the [`zonestat\(1\)`](#) man page and [Creating and Using Oracle Solaris Zones](#).

Design Considerations for Resource Management Applications in Oracle Solaris Zones

This chapter provides a brief overview of Oracle Solaris Zones technology and discusses potential problems that may be encountered by developers who are writing resource management applications.

Oracle Solaris Zones Overview

A *zone* is a virtualized operating system environment that is created within a single instance of the Oracle Solaris operating system. Oracle Solaris Zones are a partitioning technology that provides an isolated, secure environment for applications. When you create a zone, you produce an application execution environment in which processes are isolated from the rest of the system. This isolation prevents a process that is running in one zone from monitoring or affecting processes that are running in other zones. Even a process running with root credentials cannot view or affect activity in other zones. A zone also provides an abstract layer that separates applications from the physical attributes of the system on which the zone is deployed. Examples of these attributes include physical device paths and network interface names. The default non-global zone brand in the Oracle Solaris 11.3 release is the `solaris` zone.

By default, all systems have a *global zone*. The global zone has a global view of the Oracle Solaris environment that is similar to the superuser (root) model. All other zones are referred to as *non-global zones*. A non-global zone is analogous to an unprivileged user in the superuser model. Processes in non-global zones can control only the processes and files within that zone. Typically, system administration work is mainly performed in the global zone. In rare cases where a system administrator needs to be isolated, privileged applications can be used in a non-global zone. In general, though, resource management activities take place in the global zone.

For additional isolation, `solaris` zones with a read-only root can be configured. See [Chapter 11, “Configuring and Administering Immutable Zones”](#) in *Creating and Using Oracle Solaris Zones*.

For more information on `solaris` zones, see *Creating and Using Oracle Solaris Zones*.

IP Networking in Oracle Solaris Zones

IP networking in a zone can be configured in two different ways, depending on whether the non-global zone is given its own exclusive IP instance or shares the IP layer configuration and state with the global zone. By default, zones are created with the exclusive-IP type. Through the `zonecfg anet` resource, a virtual network (VNIC) is automatically included in the zone configuration if networking configuration is not specified.

Exclusive-IP zones are assigned zero or more VNIC interface names, and for those network interfaces they can send and receive any packets, snoop, and change the IP configuration, including IP addresses and the routing table. Note that those changes do not affect any of the other IP instances on the system.

For complete information on the `zonecfg` command and networking in zones, see [Oracle Solaris Zones Configuration Resources](#).

About Applications in Oracle Solaris Zones

All applications are fully functional in the global zone, as they would be in a conventional Oracle Solaris operating system. Most applications should run without problem in a non-global environment as long as the application does not need any privileges. If an application does require privileges, then the developer needs to take a close look at which privileges are needed and how a particular privilege is used. If a privilege is required, then a system administrator can assign the needed privilege to the zone. See [Introduction to Oracle Solaris Zones](#).

General Considerations When Writing Applications for Non-Global Zones

The known situations that a developer should investigate are as follows:

- System calls that change the system time require the PRIV_SYS_TIME privilege. These system calls include [adjtime\(2\)](#), [ntp_adjtime\(2\)](#), and [stime\(2\)](#).
- System calls that need to operate on files that have the sticky bit set require the PRIV_SYS_CONFIG privilege. These system calls include [chmod\(2\)](#), [creat\(2\)](#), and [open\(2\)](#).
- The [ioctl\(2\)](#) system call requires the PRIV_SYS_NET_CONFIG privilege to be able to unlock an anchor on a STREAMS module.
- The [link\(2\)](#) and [unlink\(2\)](#) system calls require the PRIV_SYS_LINKDIR privilege to create a link or unlink a directory in a non-global zone. Applications that install or configure software or that create temporary directories could be affected by this limitation.
- The PRIV_PROC_LOCK_MEMORY privilege is required for the [mlock\(3C\)](#), [munlock\(3C\)](#), [mlockall\(3C\)](#), [munlockall\(3C\)](#), and [plock\(3C\)](#) functions and the MC_LOCK, MC_LOCKAS, MC_UNLOCK, and MC_UNLOCKAS flags for the [memcntl\(2\)](#) system. This privilege is a default privilege in a non-global zone. See “Privileges in a Non-Global Zone” in *Creating and Using Oracle Solaris Zones* for more information.
- The [mknod\(2\)](#) system call requires the PRIV_SYS_DEVICES privilege to create a block (S_IFBLK) or character (S_IFCHR) special file. This limitation affects applications that need to create device nodes on the fly.
- The IPC_SET flag in the [msgctl\(2\)](#) system call requires the PRIV_SYS_IPC_CONFIG privilege to increase the number of message queue bytes. This limitation affects any applications that need to resize the message queue dynamically.
- The [nice\(2\)](#) system call requires the PRIV_PROC_PRIOCNTL privilege to change the priority of a process. This privilege is available by default in a non-global zone. Another way to change the priority is to bind the non-global zone in which the application is running to a resource pool, although scheduling processes in that zone is ultimately decided by the Fair Share Scheduler.
- The P_ONLINE, P_OFFLINE, P_NOINTR, P_FAULTED, P_SPARE, and PZ-FORCED flags in the [p_online\(2\)](#) system call require the PRIV_SYS_RES_CONFIG privilege to return or change process operational status. This limitation affects applications that need to enable or disable CPUs.
- The PC_SETPARMS and PC_SETXPARMS flags in the [priocntl\(2\)](#) system call requires the PRIV_PROC_PRIOCNTL privilege to change the scheduling parameters of a lightweight process (LWP).
- System calls that need to manage processor sets (psets), including binding LWPs to psets and setting pset attributes require the PRIV_SYS_RES_CONFIG privilege. This limitation

affects the following system calls: [pset_assign\(2\)](#), [pset_bind\(2\)](#), [pset_create\(2\)](#), [pset_destroy\(2\)](#), and [pset_setattr\(2\)](#).

- The SHM_LOCK and SHM_UNLOCK flags in the [shmctl\(2\)](#) system call require the PRIV_PROC_LOCK_MEMORY privilege to share memory control operations. If the application is locking memory for performance purposes, using the intimate shared memory (ISM) feature provides a potential workaround.
- The [swapctl\(2\)](#) system call requires the PRIV_SYS_CONFIG privilege to add or remove swapping resources. This limitation affects installation and configuration software.
- The [uadmin\(2\)](#) system call requires the PRIV_SYS_CONFIG privilege to use the A_REMOUNT, A_FREEZE, A_DUMP, and AD_IBOOT commands. This limitation affects applications that need to force crash dumps under certain circumstances.
- The [clock-settime\(3C\)](#) function requires the PRIV_SYS_TIME privilege to set the CLOCK_REALTIME and CLOCK_HIRES clocks.
- The [cpc_bind_cpu\(3CPC\)](#) function requires the PRIV_CPC_CPU privilege to bind request sets to hardware counters. As a workaround, the [cpc_bind_curlwp\(3CPC\)](#) function can be used to monitor CPU counters for the LWP in question.
- The [pthread_attr_setschedparam\(3C\)](#) function requires the PRIV_PROC_PRIOCNTRL privilege to change the underlying scheduling policy and parameters for a thread.
- The [timer_create\(3C\)](#) function requires the PRIV_PROC_CLOCK_HIGHRES privilege to create a timer using the high-resolution system clock.
- The APIs that are provided by the following list of libraries are not supported in a non-global zone. The shared objects are present in the zone's `/usr/lib` directory, so no link time errors occur if your code includes references to these libraries. You can inspect your make files to determine if your application has explicit bindings to any of these libraries and use [pmap\(1\)](#) while the application is executing to verify that none of these libraries are dynamically loaded.
 - [libdevinfo\(3LIB\)](#)
 - [libcfgadm\(3LIB\)](#)
 - [libpool\(3LIB\)](#)
 - [libsysevent\(3LIB\)](#)
- Zones have a restricted set of devices, consisting primarily of pseudo devices that form part of the Oracle Solaris programming API. These pseudo devices include `/dev/null`, `/dev/zero`, `/dev/poll`, `/dev/random`, `/dev/tcp`, and so on. Physical devices are not directly accessible from within a zone unless the device has been configured by a system administrator. Since devices, in general, are shared resources in a system, to make devices available in a zone requires some restrictions so system security will not be compromised, as follows:
 - The `/dev` name space consists of symbolic links, that is, logical paths, to the physical paths in `/devices`. The `/devices` name space, which is available only in the global

zone, reflects the current state of attached device instances that have been created by the driver. Only the logical path `/dev` is visible in a non-global zone.

- Processes within a non-global zone cannot create new device nodes. For example, `mknod(2)` cannot create special files in a non-global zone. The `creat(2)`, `link(2)`, `mknod(2)`, `rename(2)`, `symlink(2)`, and `unlink(2)` system calls fail with `EACCES` if a file in `/dev` is specified. You can create a symbolic link to an entry in `/dev`, but that link cannot be created in `/dev`.
- Devices that expose system data are only available in the global zone. Examples of such devices include `dtrace(7D)`, `kmem(7D)`, `kldb(7d)`, `ksyms(7D)`, `lockstat(7D)`, and `trapstat(1M)`.
- The `/dev` name space consists of device nodes made up of a default, "safe" set of drivers as well as device nodes that have been specified for the zone by the `zonecfg(1M)` command.

Specific Considerations for Oracle Solaris 10 Non-Global Zones

`LIFC_UNDER_IPMP` is not available in Oracle Solaris 10, so applications supported on this release do not use `LIFC_UNDER_IPMP`. Thus, applications issuing a `SIOCGLIFCONF` request inside an Oracle Solaris 10 zone do not see underlying interfaces, but instead see only the IPMP meta-interfaces. The `ifconfig` command used with the `-a` option shows the underlying interfaces because a `solaris10` zone uses the Oracle Solaris 11 version of `ifconfig`, which passes the special `LIFC_UNDER_IPMP`. See `if_tcp(7P)` for details.

Specific Considerations for Shared-IP Non-Global Zones

For non-global zones that are configured to use the shared-IP instance, the following restrictions apply.

- The `socket(3SOCKET)` function requires the `PRIV_NET_RAWACCESS` privilege to create a raw socket with the protocol set to `IPPROTO_RAW` or `IPPROTO_IGMP`. This limitation affects applications that use raw sockets or need to create or inspect TCP/IP headers.
- The `t_open(3NSL)` function requires the `PRIV_NET_RAWACCESS` privilege to establish a transport endpoint. This limitation affects applications that use the `/dev/rawip` device to implement network protocols as well as applications that operate on TCP/IP headers.

- No NIC devices that support the DLPI programming interface are accessible in a shared-IP non-global zone.
- Each non-global shared-IP zone has its own logical network and loopback interface. Bindings between upper layer streams and logical interfaces are restricted such that a stream may only establish bindings to logical interfaces in the same zone. Likewise, packets from a logical interface can only be passed to upper layer streams in the same zone as the logical interface. Bindings to the loopback address are kept within a zone with one exception: When a stream in one zone attempts to access the IP address of an interface in another zone. While applications within a zone can bind to privileged network ports, they have no control over the network configuration, including IP addresses and the routing table.

Note that these restrictions do not apply to exclusive-IP zones.

Packaging Considerations in solaris Zones

Using a zone package variant, the various components within a package are specifically tagged to only be installed in either a global zone (`global`) or a non-global zone (`nonglobal`). A given package can contain a file that is tagged so that it will not be installed into a non-global zone.

API for Zones Monitoring Statistics

`libzonestat.so.1` is a public API used by the `zonestat` command to retrieve and compute zone-related resource utilization information, with sorting and filtering options available. The `zonestat` library reports system wide and per-zone utilization of physical memory, virtual memory, and CPU resources. The `zonestat` command is documented in the [zonestat\(1\)](#) man page.

`libzonestat` computes commonly needed values, such as differences between two samples, and percentage used quantities. These statistics eliminate the need for consumers to do complex calculations. In addition to usage of physical resources, `libzonestat` also reports resource usage relative to each zone's configured resource limits.

The basic usage of `libzonestat` is as follows:

```
zs_ctl_t zsctl;
zs_usage_t usage;

/* open the statistics facility */
zsctl = zs_open();
```

```

    for (;;) {
        /* read the current usage */
        usage = zs_usage_read(ctl);

        ... Interrogate the usage object for desired information ...
        ...

        if (quit)
            break;

        sleep(some_interval);
    }
    zs_close(zsctl);

```

The following man pages, ordered to facilitate usage, describe the library interfaces:

- [zs_open\(3ZONESTAT\)](#)
- [zs_usage\(3ZONESTAT\)](#)
- [libzonestat\(3LIB\)](#)
- [zs_resource\(3ZONESTAT\)](#)
- [zs_zone\(3ZONESTAT\)](#)
- [zs_pset\(3ZONESTAT\)](#)
- [zs_pset_zone\(3ZONESTAT\)](#)
- [zs_property\(3ZONESTAT\)](#)

Monitoring Zone File System Activity

The `fsstat` utility can be used to report file operations statistics for non-global zones.

The global zone can see the `kstats` of all zones on the system. Non-global zones only see the `kstats` associated with the zone in which the utility is run. A non-global zone cannot monitor file system activity in other zones.

The `-z` option is used to report on file system activity per zone. Multiple `-z` options can be used to monitor activity in selected zones. The option has no effect if only used to monitor mountpoints and not `fstypes`.

The `-Z` option is used to report file system activity in all zones on the system. This option has no effect if used with `-z` option. The option has no effect if only used to monitor mountpoints and not `fstypes`.

The `-A` option is used to report on aggregate file system activity for the specified `fstypes` across all zones. This is the default behavior if neither `-z` or the `-Z` option is used. The `-A` option has no effect if only used to monitor mountpoints and not `fstypes`.

When used with either the `-z` or the `-Z` option, the `-A` option displays the aggregate for the specified `fstypes` across all zones on a separate line. The option has no effect if only used to monitor mountpoints and not `fstypes`.

See the [fsstat\(1M\)](#) man page for more information.

Oracle Solaris 10 Zones

Oracle Solaris 10 Zones are `solaris10` branded zones that host x86 and SPARC Solaris 10 9/10 (or later released Oracle Solaris 10 update) user environments running on the Oracle Solaris 11 kernel. Note that it is possible to use an earlier Oracle Solaris 10 release if you first install the kernel patch 142909-17 (SPARC) or 142910-17 (x86/x64), or later version, on the original system.

For more information about `solaris10` branded zones, see the following guides:

- [Creating and Using Oracle Solaris 10 Zones](#).
- [System Administration Guide: Oracle Solaris Containers-Resource Management and Oracle Solaris Zones](#) (this is the Oracle Solaris 10 version of the guide).
- [Solaris Containers: Resource Management and Solaris Zones Developer's Guide](#) (this is the Oracle Solaris 10 version of the guide).

For information about SVR4 packaging and patching used in `solaris10` and native zones, see "Chapter 25, About Packages on an Solaris System With Zones Installed (Overview)" and "Chapter 26, Adding and Removing Packages and Patches on a Solaris System With Zones Installed (Tasks)" in the Oracle Solaris 10 guide [System Administration Guide: Oracle Solaris Containers-Resource Management and Oracle Solaris Zones](#).

Oracle Solaris Kernel Zones

The Oracle Solaris Kernel Zones feature provides a full kernel and user environment within a zone. Kernel zones also increase kernel separation between the host system and the zone. The brand name is `solaris-kz`. Kernel zones are managed by using the existing tools. You can fully update and modify the zone's installed packages, including the kernel version, without

being limited to the packages installed in the global zone. You can run `solaris` zones within the kernel zone to produce hierarchical (nested) zones. Kernel zones support suspend and resume. The Oracle Solaris Kernel Zone can run a Support Repository Update (SRU) or kernel version that is different from that of the host system.

For information about Oracle Solaris Kernel Zones, see [Creating and Using Oracle Solaris Kernel Zones](#) and [Chapter 1, “Configuration Resources for Non-Global Zones”](#) in [Oracle Solaris Zones Configuration Resources](#).

Project Configuration Examples

This chapter shows example configurations for the `/etc/project` file.

- “Configure Resource Controls” on page 100
- “Configure Resource Pools” on page 100
- “Configure FSS `project.cpu-shares` for a Project” on page 100
- “Configure Five Applications with Different Characteristics” on page 101

To set resource controls on a zone by using the `zonectfg` command, see [Chapter 1](#), “Configuration Resources for Non-Global Zones” in *Oracle Solaris Zones Configuration Resources* and “How to Configure the Zone” in *Creating and Using Oracle Solaris Zones*.

`/etc/project` Project File

The project file is a local source of project information. The project file can be used in conjunction with other project sources, including the NIS maps `project.byname` and `project.bynumber` and the LDAP database `project`. Programs use the `getproject(3PROJECT)` routines to access this information.

Define Two Projects

`/etc/project` defines two projects: `database` and `appserver`. The *user* defaults are `user.database` and `user.appserver`. The *admin* default can switch between `user.database` or `user.appserver`.

```
hostname# cat /etc/project
.
.
.
user.database:2001:Database backend:admin::
```

```
user.appserver:2002:Application Server frontend:admin::
```

```
.  
.
```

Configure Resource Controls

The /etc/project file shows the resource controls for the application.

```
hostname# cat /etc/project
```

```
.  
.  
.  
development:2003:Developers:::task.ax-lwps=(privileged,10,deny);  
process.max-addressspace=(privileged,209715200,deny)  
.  
.
```

Configure Resource Pools

The /etc/project file shows the resource pools for the application.

```
hostname# cat /etc/project
```

```
.  
.  
.  
batch:2001:Batch project:::project.pool=batch_pool  
process:2002:Process control:::project.pool=process_pool  
.  
.  
.
```

Configure FSS project.cpu-shares for a Project

Set up FSS for two projects: *database* and *appserver*. The *database* project has 20 CPU shares. The *appserver* project has 10 CPU shares.

```
hostname# cat /etc/project
```

```
.  
.
```

```

.
user.database:2001:database backend:admin::project.cpu-shares=(privileged,
20,deny)
user.appserver:2002:Application Server frontend:admin::project.cpu-shares=
(privileged,10,deny)
.
.
.

```

Note - The line break in the lines that precede “20,deny” and “(privileged,” is not valid in an /etc/project file. The line breaks are shown here only to allow the example to display on a printed or displayed page. Each entry in the /etc/project file must be on a single line.

If the FSS is enabled but each user and application is not assigned to a unique project, then the users and applications will all run in the same project. By running in the same project, all compete for the same share, in a timeshare fashion. This occurs because shares are assigned to projects, not to users or processes. To take advantage of the FSS scheduling capabilities, assign each user and application to a unique project.

To configure a project, see [“Local /etc/project File Format” in *Administering Resource Management in Oracle Solaris 11.3*](#).

Configure Five Applications with Different Characteristics

The following example configures five applications with different characteristics.

TABLE 8 Target Applications and Characteristics

Application Type and Name	Characteristics
Application server, app_server.	Negative scalability beyond two CPUs. Assign a two-CPU processor set to app_server. Use TS scheduling class.
Database instance, app_db.	Heavily multithreaded. Use FSS scheduling class.
Test and development, development.	Motif based. Hosts untested code execution. Interactive scheduling class ensures user interface responsiveness. Use process.max-address-space to impose memory limitations and minimize the effects of antisocial processing.
Transaction processing engine, tp_engine.	Response time is paramount. Assign a dedicated set of at least two CPUs to ensure response latency is kept to a minimum. Use timeshare scheduling class.
Standalone database instance, geo_db.	Heavily multithreaded. Serves multiple time zones. Use FSS scheduling class.

Note - Consolidate database applications (app_db and geo_db) onto a single processor set of at least four CPUs. Use FSS scheduling class. Application app_db gets 25% of the project.cpu-shares. Application geo_db gets 75% of the project.cpu-shares.

Edit the /etc/project file. Map users to resource pools for the app_server, app_db, development, tp_engine, and geo_db project entries.

```
hostname# cat /etc/project
```

```
.  
. .  
user.app_server:2001:Production Application Server::  
    project.pool=appserver_pool  
user.app_db:2002:App Server DB::project.pool=db_pool,  
    project.cpu-shares=(privileged,1,deny)  
development:2003:Test and delopment::staff:project.pool=dev.pool,  
    process.max-addressspace=(privileged,536870912,deny)  
user.tp_engine:Transaction Engine::project.pool=tp_pool  
user.geo_db:EDI DB::project.pool=db_pool;  
    project.cpu-shares=(privileged,3,deny)
```

Note - The line break in the lines that begin with “project.pool”, “project.cpu-shares=”, “process.max-addressspace”, and “project.cpu-shares=” is not valid in a project file. The line breaks are shown here only to allow the example to display on a printed or displayed page. Each entry must be on one and only one line.

Create the pool.host script and add entries for resource pools.

```
hostname# cat pool.host
```

```
create system host  
create pset dev_pset (uint pset.max = 2)  
create pset tp_pset (uint pset.min = 2; uint pset.max = 2)  
create pset db_pset (uint pset.min = 4; uint pset.max = 6)  
create pset app_pset (uint pset.min = 1; uint pset.max = 2)  
create pool dev_pool (string pool.scheduler="IA")  
create pool appserver_pool (string pool.scheduler="TS")  
create pool db_pool (string pool.scheduler="FSS")  
create pool tp_pool (string pool.scheduler="TS")  
associate pool pool_default (pset pset_default)  
associate pool dev_pool (pset dev_pset)  
associate pool appserver_pool (pset app_pset)  
associate pool db_pool (pset db_pset)  
associate pool tp_pool (pset tp_pset)
```

Run the pool.host script and modify the configuration as specified in the pool.host file.

```
hostname# poolcfg -f pool.host
```

Read the pool.host resource pool configuration file and initialize the resource pools on the system.

```
hostname# pooladm -c
```


Index

B

brands, 96

E

ea_alloc(), 23
ea_copy_object(), 23
ea_copy_object_tree(), 24
ea_free(), 23
ea_free_item(), 23
ea_free_object(), 23
ea_get_object_tree(), 24
ea_pack_object(), 23
ea_strdup(), 23
ea_strfree(), 23
ea_unpack_object(), 23
exacct file
 display entry, 24
 display string, 25
 display system file, 26
 dump, 46
exacct object
 create record, 46
 dump, 44
 write file, 46

F

fair share scheduler
 access resource control block, 67
fsstat
 zone, 95

L

libexacct
 perl interface, 30
 perl module, 31
libzonestat API, 94

O

Oracle Solaris Kernel Zones, 96
Oracle Solaris Zones
 overview, 89

P

programming issues
 exacct files, 28
 project database, 20
 resource controls, 69
project database
 get entry, 20
 print entries, 19

R

resource controls
 display value-action pairs, 66
 global action, 54
 global flag, 54
 local action, 53
 local flag, 53
 master observer process, 64
 privilege levels, 53

- process, 59
- project, 57
- signals, 62
- task, 58
- zone, 60
- resource pools
 - get defined pools, 84
 - get number of CPUS, 83
 - get pool statistics, 85
 - overview, 71
 - pool properties, 74
 - processor sets properties, 75
 - properties, 72
 - scheduling class, 72
 - set property, 85
 - system properties, 73

S

- solaris-kz brand, 96

Z

- zone
 - application design considerations, 90
 - exclusive-IP type, 90
 - fsstat, 95
 - IP type, 90
 - libzonestat, 94
 - packaging, 94
 - resource controls, 60
 - shared-IP, 93
 - solaris brand, 89
 - solaris10 brand, 96
 - zonestat, 94
- zonestat
 - zone, 94
- zonestat utility, 69, 87