



System Administration Guide: Solaris Containers—Resource Management and Solaris Zones

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

Part No: 817-1592-10
January 2005

Copyright 2005 Sun Microsystems, Inc. 4150 Network Circle, Santa Clara, CA 95054 U.S.A. All rights reserved.

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

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

Sun, Sun Microsystems, the Sun logo, SunOS, docs.sun.com, AnswerBook, AnswerBook2, Solaris, StarOffice, J2SE, and CacheFS are trademarks or registered trademarks of Sun Microsystems, Inc. in the U.S. and other countries. All SPARC trademarks are used under license and are trademarks or registered trademarks of SPARC International, Inc. in the U.S. and other countries. Products bearing SPARC trademarks are based upon an architecture developed by Sun Microsystems, Inc.

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

U.S. Government Rights – Commercial software. Government users are subject to the Sun Microsystems, Inc. standard license agreement and applicable provisions of the FAR and its supplements.

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

Copyright 2005 Sun Microsystems, Inc. 4150 Network Circle, Santa Clara, CA 95054 U.S.A. Tous droits réservés.

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

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

Sun, Sun Microsystems, the Sun logo, SunOS, docs.sun.com, AnswerBook, AnswerBook2, Solaris, StarOffice, J2SE, et CacheFS sont des marques de fabrique ou des marques déposées de Sun Microsystems, Inc. aux Etats-Unis et dans d'autres pays. Toutes les marques SPARC sont utilisées sous licence et sont des marques de fabrique ou des marques déposées de SPARC International, Inc. aux Etats-Unis et dans d'autres pays. Les produits portant les marques SPARC sont basés sur une architecture développée par Sun Microsystems, Inc.

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

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



041208@10536



Contents

Preface 17

Part I Resource Management 21

1 Introduction to Solaris 10 Resource Manager 23

What's New in Resource Management? 23

Improvements to Project Database and Resource Control Commands 23

New Resource Controls 24

Resource Pools Changes 25

New Extended Accounting Features 25

Ability to Regulate Physical Memory Consumption by Processes 25

Interaction With Solaris Zones 25

Resource Management Overview 26

Resource Classifications 27

Resource Management Control Mechanisms 27

Resource Management Configuration 28

When to Use Resource Management 29

Server Consolidation 29

Supporting a Large or Varied User Population 29

Setting Up Resource Management (Task Map) 30

2 Projects and Tasks (Overview) 33

Project and Task Facilities 33

Project Identifiers 34

Determining a User's Default Project 34

Setting User Attributes With the <code>useradd</code> , <code>usermod</code> , and <code>passmgmt</code> Commands	35
project Database	35
PAM Subsystem	36
Naming Services Configuration	36
Local <code>/etc/project</code> File Format	36
Project Configuration for NIS	38
Project Configuration for LDAP	39
Task Identifiers	39
Commands Used With Projects and Tasks	40
3 Administering Projects and Tasks	43
Administering Projects and Tasks (Task Map)	43
Sample Commands and Command Options	44
Command Options Used With Projects and Tasks	44
Using <code>cron</code> and <code>su</code> With Projects and Tasks	46
Administering Projects	47
▼ How to Define a Project and View the Current Project	47
▼ How to Delete a Project From the <code>/etc/project</code> File	49
How to Validate the Contents of the <code>/etc/project</code> File	50
How to Obtain Project Membership Information	51
▼ How to Create a New Task	51
▼ How to Move a Running Process Into a New Task	51
Editing and Validating Project Attributes	52
How to Add Attributes and Attribute Values to Projects	52
How to Remove Attributes and Attribute Values From Projects	53
How to Substitute Attributes and Attribute Values for Projects	53
4 Extended Accounting (Overview)	55
Introduction to Extended Accounting	55
How Extended Accounting Works	56
Extensible Format	57
<code>exacct</code> Records and Format	57
Using Extended Accounting in a Zones Environment	58
Extended Accounting Configuration	58
Commands Used With Extended Accounting	59
Perl Interface to <code>libexacct</code>	59

5	Administering Extended Accounting (Tasks)	63
	Administering the Extended Accounting Facility (Task Map)	63
	Using Extended Accounting Functionality	64
	▼ How to Activate Extended Accounting for Processes, Tasks, and Flows	64
	How to Activate Extended Accounting With a Startup Script	64
	How to Display Extended Accounting Status	65
	How to View Available Accounting Resources	65
	▼ How to Deactivate Process, Task, and Flow Accounting	66
	Using the Perl Interface to <code>libexacct</code>	67
	How to Recursively Print the Contents of an <code>exacct</code> Object	67
	How to Create a New Group Record and Write It to a File	68
	How to Print the Contents of an <code>exacct</code> File	69
	Example Output From <code>Sun::Solaris::Exacct::Object->dump()</code>	70
6	Resource Controls (Overview)	71
	Resource Controls Concepts	71
	Resource Limits and Resource Controls	72
	Interprocess Communication and Resource Controls	72
	Resource Control Constraint Mechanisms	72
	Project Attribute Mechanisms	73
	Configuring Resource Controls and Attributes	73
	Available Resource Controls	74
	Units Support	77
	Resource Control Values and Privilege Levels	78
	Global and Local Actions on Resource Control Values	79
	Resource Control Flags and Properties	81
	Resource Control Enforcement	82
	Global Monitoring of Resource Control Events	82
	Applying Resource Controls	83
	Temporarily Updating Resource Control Values on a Running System	83
	Updating Logging Status	83
	Updating Resource Controls	83
	Commands Used With Resource Controls	84
7	Administering Resource Controls (Tasks)	85
	Administering Resource Controls (Task Map)	85
	Setting Resource Controls	86

▼ How to Set the Maximum Number of LWPs for Each Task in a Project	86
▼ How to Set Multiple Controls on a Project	87
Using the <code>prctl</code> Command	89
How to Use <code>prctl</code> to Display Information	89
How to Use <code>prctl</code> to Temporarily Change a Value	89
How to Use <code>prctl</code> to Lower a Resource Control Value	90
▼ How to Use <code>prctl</code> to Display, Replace, and Verify the Value of a Control on a Project	90
Using <code>rctladm</code>	91
How to Use <code>rctladm</code>	91
Using <code>ipcs</code>	92
How to Use <code>ipcs</code>	92
Capacity Warnings	92
▼ How to Determine Whether a Web Server Is Allocated Enough CPU Capacity	93
8 Fair Share Scheduler (Overview)	95
Introduction to the Scheduler	95
CPU Share Definition	96
CPU Shares and Process State	97
CPU Share Versus Utilization	97
CPU Share Examples	98
Example 1: Two CPU-Bound Processes in Each Project	98
Example 2: No Competition Between Projects	99
Example 3: One Project Unable to Run	99
FSS Configuration	100
Projects and Users	100
CPU Shares Configuration	100
FSS and Processor Sets	102
FSS and Processor Sets Examples	103
Combining FSS With Other Scheduling Classes	104
Setting the Scheduling Class for the System	105
Scheduling Class in a Zones Environment	105
Commands Used With FSS	106
9 Administering the Fair Share Scheduler (Tasks)	107
Administering the Fair Share Scheduler (Task Map)	107
Monitoring the FSS	108

	▼ How to Monitor System CPU Usage by Projects	108
	▼ How to Monitor CPU Usage by Projects in Processor Sets	108
	Configuring the FSS	109
	▼ How to Set the Scheduler Class	109
	▼ How to Manually Move Processes From the TS Class Into the FSS Class	109
	▼ How to Manually Move Processes From All User Classes Into the FSS Class	110
	▼ How to Manually Move a Project's Processes Into the FSS Class	110
	How to Tune Scheduler Parameters	110
10	Physical Memory Control Using the Resource Capping Daemon (Overview)	113
	Introduction to the Resource Capping Daemon	113
	How Resource Capping Works	114
	Attribute to Limit Physical Memory Usage	115
	rcapd Configuration	115
	Memory Cap Enforcement Threshold	116
	Determining Cap Values	116
	rcapd Operation Intervals	118
	Monitoring Resource Utilization With rcapstat	119
	Commands Used With rcapd	121
11	Administering the Resource Capping Daemon (Tasks)	123
	Configuring and Using the Resource Capping Daemon (Task Map)	123
	Administering the Resource Capping Daemon With rcapadm	124
	▼ How to Set the Memory Cap Enforcement Threshold	124
	▼ How to Set Operation Intervals	125
	▼ How to Enable Resource Capping	125
	▼ How to Disable Resource Capping	126
	Producing Reports With rcapstat	126
	Reporting Cap and Project Information	126
	Monitoring the RSS of a Project	127
	Determining the Working Set Size of a Project	128
	Reporting Memory Utilization and the Memory Cap Enforcement Threshold	130
12	Dynamic Resource Pools (Overview)	133
	Introduction to Resource Pools	133

project.pool Attribute	135
Resource Pools Used in Zones	135
When to Use Pools	136
Resource Pools Framework	137
/etc/pooladm.conf Contents	138
Pools Properties	138
Implementing Pools on a System	139
SPARC: Dynamic Reconfiguration Operations and Resource Pools	140
Creating Pools Configurations	140
Directly Manipulating the Dynamic Configuration	141
poold Overview	141
Stopping poold	142
Reconfiguring poold	142
Configuration Constraints and Objectives	142
Configuration Constraints	143
Configuration Objectives	143
poold Properties	146
poold Features That Can Be Configured	147
poold Monitoring Interval	147
poold Logging Information	148
Logging Location	150
Log Management With logadm	150
How Dynamic Resource Allocation Works	150
About Available Resources	150
Determining Available Resources	151
Identifying a Resource Shortage	152
Determining Resource Utilization	152
Identifying Control Violations	152
Determining Appropriate Remedial Action	153
Using poolstat to Monitor the Pools Facility and Resource Utilization	153
poolstat Output	154
Tuning poolstat Operation Intervals	154
Commands Used With the Resource Pools Facility	155
13 Administering Dynamic Resource Pools (Tasks)	157
Administering Dynamic Resource Pools (Task Map)	157
Enabling and Disabling the Pools Facility	159
▼ How to Enable Pools	159

	▼ How to Disable Pools	159
	Configuring Pools	160
	▼ How to Create a Static Configuration	160
	▼ How to Modify a Configuration	161
	▼ How to Associate a Pool With a Scheduling Class	163
	▼ How to Define Configuration Objectives	165
	▼ How to Set the <code>poold</code> Logging Level	168
	▼ How to Use Command Files With <code>poolcfg</code>	168
	Transferring Resources	169
	Activating and Removing Pool Configurations	169
	▼ How to Activate a Pools Configuration	170
	▼ How to Validate a Configuration Before Committing the Configuration	170
	▼ How to Remove a Pools Configuration	170
	Binding to a Pool	171
	▼ How to Bind Processes to a Pool	171
	▼ How to Bind Tasks or Projects to a Pool	172
	▼ How to Use <code>project</code> Attributes to Bind New Processes to a Pool	172
	▼ How to Use <code>project</code> Attributes to Bind a Process to a Different Pool	173
	Using <code>poolstat</code> to Report Statistics for Pool-Related Resources	173
	Displaying Default <code>poolstat</code> Output	173
	Producing Multiple Reports at Specific Intervals	173
	Reporting Resource Set Statistics	174
14	Resource Management Configuration Example	175
	Configuration to Be Consolidated	175
	Consolidation Configuration	176
	Creating the Configuration	176
	Viewing the Configuration	178
15	Resource Control Functionality in the Solaris Management Console	183
	Using the Console (Task Map)	184
	Console Overview	184
	Management Scope	184
	Performance Tool	185
	▼ How to Access the Performance Tool	185
	Monitoring by System	186
	Monitoring by Project or User Name	186

Resource Controls Tab	188
▼ How to Access the Resource Controls Tab	189
Resource Controls You Can Set	190
Setting Values	191
Console References	191

Part II Zones 193

16 Introduction to Solaris Zones	195
Zones Overview	195
When to Use Zones	196
How Zones Work	198
Summary of Zone Features	199
How Non-Global Zones Are Administered	200
How Non-Global Zones Are Created	200
Non-Global Zone State Model	201
Non-Global Zone Characteristics	202
Using Resource Management Features With Non-Global Zones	203
Features Provided by Non-Global Zones	203
Setting Up Zones on Your System (Task Map)	204
17 Non-Global Zone Configuration (Overview)	207
Pre-Installation Configuration Process	207
Zone Components	208
Zone Name and Path	208
Zone Interfaces	208
File Systems Mounted in Zones	208
Configured Devices in Zones	209
Zone-Wide Resource Controls	209
Including a Comment for a Zone	209
Using the zonecfg Command	209
zonecfg Modes	210
zonecfg Interactive Mode	210
zonecfg Command-File Mode	212
Zone Configuration Data	212
Resource Types	213
Resource Type Properties	214

Tecla Command-Line Editing Library 217

- 18 Planning and Configuring Non-Global Zones (Tasks) 219**
 - Planning and Configuring a Non-Global Zone (Task Map) 219
 - Evaluating the Current System Setup 221
 - Disk Space Requirements 221
 - Restricting Zone Size 222
 - Determine the Zone Host Name and Obtain the Network Address 223
 - Zone Host Name 223
 - Zone Network Address 223
 - File System Configuration 224
 - Creating, Revising, and Deleting Non-Global Zone Configurations (Task Map) 225
 - Configuring, Verifying, and Committing a Zone 226
 - ▼ How to Configure the Zone 226
 - Script to Configure Multiple Zones 230
 - Using the `zonecfg` Command to Modify a Zone Configuration 232
 - ▼ How to Modify a Resource Type in a Zone Configuration 232
 - ▼ How to Add a Dedicated Device to a Zone 233
 - Using the `zonecfg` Command to Revert or Remove a Zone Configuration 234
 - ▼ How to Revert a Zone Configuration 234
 - ▼ How to Delete a Zone Configuration 236

- 19 About Installing, Halting, and Uninstalling Non-Global Zones (Overview) 237**
 - Zone Installation Concepts 237
 - Zone Construction 238
 - The `zoneadmd` Daemon 239
 - The `zsched` Zone Scheduler 240
 - Zone Application Environment 240
 - About Halting, Rebooting, and Uninstalling Zones 240
 - Halting a Zone 241
 - Rebooting a Zone 241
 - Zone `autoboot` 241
 - Uninstalling a Zone 241

- 20 Installing, Booting, Halting, and Uninstalling Non-Global Zones (Tasks) 243**
 - Zone Installation (Task Map) 243
 - Installing and Booting Zones 244

▼ How to Verify a Configured Zone Before It Is Installed (Optional)	244
▼ How to Install a Configured Zone	245
▼ How to Transition the Installed Zone to the Ready State (Optional)	246
▼ How to Boot a Zone	246
Where to Go From Here	248
Halting, Rebooting, Uninstalling, and Deleting Non-Global Zones (Task Map)	248
Halting, Rebooting, and Uninstalling Zones	249
▼ How to Halt a Zone	249
▼ How to Reboot a Zone	250
▼ How to Uninstall a Zone	250
Deleting a Non-Global Zone From the System	251
▼ How to Remove a Non-Global Zone	251
21 Non-Global Zone Login (Overview)	253
zlogin Command	253
Internal Zone Configuration	254
Non-Global Zone Login Methods	255
Zone Console Login	255
User Login Methods	255
Failsafe Mode	255
Remote Login	256
Interactive and Non-Interactive Modes	256
Interactive Mode	256
Non-Interactive Mode	256
22 Logging In to Non-Global Zones (Tasks)	257
Initial Zone Boot and Zone Login Procedures (Task Map)	257
Performing the Initial Internal Zone Configuration	258
▼ How to Log In to the Zone Console to Perform the Internal Zone Configuration	258
▼ How to Use an <code>/etc/sysidcfg</code> File to Perform the Initial Zone Configuration	260
Logging In to a Zone	261
▼ How to Log In to the Zone Console	261
▼ How to Use Interactive Mode to Access a Zone	262
▼ How to Use Non-Interactive Mode to Access a Zone	263
▼ How to Use Failsafe Mode to Enter a Zone	263
▼ How to Use <code>zlogin</code> to Shut Down a Zone	264

Printing the Name of the Current Zone 264

23 About Adding and Removing Packages and Patches in a Solaris Zones Environment (Overview) 265

Packaging and Patch Tools Overview 265

About Packages and Zones 267

Keeping Zones in Sync 267

 Package Operations Possible in the Global Zone 268

 Package Operations Possible in a Non-Global Zone 268

About Adding Packages in Zones 269

 Using `pkgadd` in the Global Zone 269

 Using `pkgadd` in a Non-Global Zone 271

About Removing Packages in Zones 272

 Using `pkgrm` in the Global Zone 272

 Using `pkgrm` in a Non-Global Zone 274

Package Parameter Information 274

 SUNW_PKG_ALLZONES Package Parameter 274

 SUNW_PKG_HOLLOW Package Parameter 276

Package Information Query 278

 Using `pkginfo` in the Global Zone 278

 Using `pkginfo` in a Non-Global Zone 278

About Adding Patches in Zones 279

Applying Patches in a Zones Environment 279

 Using `patchadd` in the Global Zone 280

 Using `patchadd` in a Non-Global Zone 280

Removing Patches in a Zones Environment 280

 Using `patchrm` in the Global Zone 281

 Using `patchrm` in a Non-Global Zone 281

PatchPro Support 281

Product Database 281

24 Adding and Removing Packages and Patches in a Solaris Zones Environment (Tasks) 283

Adding and Removing Packages and Patches in a Zones Environment (Task Map) 283

Adding a Package in a Zones Environment 284

 ▼ How to Add a Package to the Global Zone Only 284

 ▼ How to Add a Package to the Global Zone and All Non-Global Zones 285

- ▼ How to Add a Package That Is Installed in the Global Zone to All Non-Global Zones 285
 - ▼ How to Add a Package To a Specified Non-Global Zone Only 286
- Checking Package Information in a Zones Environment 286
 - ▼ How to Check Package Information in the Global Zone Only 286
 - ▼ How to Check Package Information in All Non-Global Zones Only 287
 - ▼ How to Check Package Information in a Specified Non-Global Zone Only 287
- Removing a Package in a Zones Environment 287
 - ▼ How to Remove a Package From the Global Zone and All Non-Global Zones 287
 - ▼ How to Remove a Package From All Non-Global Zones 288
 - ▼ How to Remove a Package From a Specified Non-Global Zone Only 288
- Applying a Patch in a Zones Environment 289
 - ▼ How to Apply a Patch to the Global Zone Only 289
 - ▼ How to Apply a Patch to the Global Zone and All Non-Global Zones 289
 - ▼ How to Apply a Patch to a Specified Non-Global Zone Only 289
- Removing a Patch in a Zones Environment 290
 - ▼ How to Remove a Patch From a Specified Non-Global Zone Only 290
- Checking Package Parameter Settings in a Zones Environment 291
 - ▼ (Optional) How to Check the Setting of a Package Already Installed on the System 291
 - ▼ (Optional) How to Check the Setting of a Package in Software on a CD-ROM 291

25 Solaris Zones Administration (Overview) 293

- Global Zone Visibility and Access 293
- Process ID Visibility in Zones 294
- System Observability in Zones 294
- Non-Global Zone Node Name 295
- File Systems and Non-Global Zones 295
 - The `-o nosuid` Option 295
 - Mounting File Systems in Zones 296
 - Security Restrictions and File System Behavior 298
 - Non-Global Zones as NFS Clients 299
 - Use of `mknod` Prohibited in a Zone 300
 - Traversing File Systems 300
 - Restriction on Accessing A Non-Global Zone From the Global Zone 300
- Networking in Non-Global Zones 301

Zone Partitioning	301
Network Interfaces	301
IP Traffic Between Zones on the Same Machine	302
IP Network Multipathing in Zones	302
Device Use in Non-Global Zones	303
/dev and the /devices Namespace	303
Exclusive-Use Devices	304
Device Driver Administration	304
Utilities That Do Not Work or Are Modified in Non-Global Zones	304
Running Applications in Non-Global Zones	305
Resource Controls Used in Non-Global Zones	305
Fair Share Scheduler in a Zones Environment	306
FSS Share Division in a Non-Global Zone	306
Share Balance Between Zones	307
Extended Accounting in a Zones Environment	307
Privileges in a Non-Global Zone	307
Using IP Security Architecture in Zones	308
Using Solaris Auditing in Zones	308
Configuring Audit in the Global Zone	308
Configuring User Audit Characteristics in a Non-Global Zone	309
Providing Audit Records for a Specific Non-Global Zone	309
Core Files in Zones	310
Commands Used in the Solaris Zones Environment	310
26 Solaris Zones Administration (Tasks)	315
Using the ppriv Utility	315
▼ How to List the Non-Global Zone's Privilege Set	315
▼ How to List a Non-Global Zone's Privilege Set With Verbose Output	316
Mounting File Systems in Running Non-Global Zones	319
▼ How to Import Raw and Block Devices by Using zonecfg	319
▼ How to Mount the File System Manually	320
▼ How to Place a File System in /etc/vfstab to Be Mounted When the Zone Boots	321
▼ How to Mount a File System From the Global Zone Into a Non-Global Zone	322
Using IP Network Multipathing in a Zones Environment	322
▼ How to Extend IP Network Multipathing Functionality to Non-Global Zones	322

Using the Fair Share Scheduler in a Zones Environment	323
How to Set FSS Shares in the Global Zone	323
▼ How to Balance CPU Usage Between the Global Zone and Non-Global Zones	324
Using Rights Profiles in Zone Administration	324
▼ How to Assign the Zone Management Profile	324
Example—Using Profile Shells With Zone Commands	325
Glossary	327
Index	331

Preface

System Administration Guide: Solaris Containers—Resource Management and Solaris Zones is part of a multivolume set that covers a significant part of the Solaris™ Operating System administration information. This book assumes that you have already installed the operating system and set up any networking software that you plan to use.

Note – This Solaris release supports systems that use the SPARC® and x86 families of processor architectures: UltraSPARC®, SPARC64, AMD64, Pentium, and Xeon EM64T. The supported systems appear in the *Solaris 10 Hardware Compatibility List* at <http://www.sun.com/bigadmin/hcl>. This document cites any implementation differences between the platform types.

About Solaris Containers

A Solaris Container is a complete runtime environment for applications. Solaris 10 Resource Manager and Solaris Zones software partitioning technology are both parts of the container. These components address different qualities the container can deliver and work together to create a complete container. The zones portion of the container provides a virtual mapping from the application to the platform resources. Zones allow application components to be isolated from one another even though the zones share a single instance of the Solaris Operating System. Resource management features permit you to allocate the quantity of resources that a workload receives.

The container establishes boundaries for resource consumption, such as CPU. These boundaries can be expanded to adapt to changing processing requirements of the application running in the container.

Who Should Use This Book

This book is intended for anyone responsible for administering one or more systems that run the Solaris 10 release. To use this book, you should have at least one to two years of UNIX[®] system administration experience.

How the System Administration Volumes Are Organized

Here is a list of the topics that are covered by the volumes of the System Administration Guides.

Book Title	Topics
<i>System Administration Guide: Basic Administration</i>	User accounts and groups, server and client support, shutting down and booting a system, managing services, and managing software (packages and patches)
<i>System Administration Guide: Advanced Administration</i>	Printing services, terminals and modems, system resources (disk quotas, accounting, and crontabs), system processes, and troubleshooting Solaris software problems
<i>System Administration Guide: Devices and File Systems</i>	Removable media, disks and devices, file systems, and backing up and restoring data
<i>System Administration Guide: IP Services</i>	TCP/IP network administration, IPv4 and IPv6 address administration, DHCP, IPsec, IKE, Solaris IP filter, Mobile IP, IP network multipathing (IPMP), and IPQoS
<i>System Administration Guide: Naming and Directory Services (DNS, NIS, and LDAP)</i>	DNS, NIS, and LDAP naming and directory services, including transitioning from NIS to LDAP and transitioning from NIS+ to LDAP
<i>System Administration Guide: Naming and Directory Services (NIS+)</i>	NIS+ naming and directory services
<i>System Administration Guide: Network Services</i>	Web cache servers, time-related services, network file systems (NFS and Autofs), mail, SLP, and PPP

Book Title	Topics
<i>System Administration Guide: Security Services</i>	Auditing, device management, file security, BART, Kerberos services, PAM, Solaris cryptographic framework, privileges, RBAC, SASL, and Solaris Secure Shell
<i>System Administration Guide: Solaris Containers-Resource Management and Solaris Zones</i>	Resource management topics projects and tasks, extended accounting, resource controls, fair share scheduler (FSS), physical memory control using the resource capping daemon (rcapd), and dynamic resource pools; virtualization using Solaris Zones software partitioning technology

Accessing Sun Documentation Online

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

Ordering Sun Documentation

Sun Microsystems offers select product documentation in print. For a list of documents and how to order them, see "Buy printed documentation" at <http://docs.sun.com>.

Typographic Conventions

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

TABLE P-1 Typographic Conventions

Typeface or Symbol	Meaning	Example
AaBbCc123	The names of commands, files, and directories, and onscreen computer output	Edit your <code>.login</code> file. Use <code>ls -a</code> to list all files. <code>machine_name%</code> you have mail.
AaBbCc123	What you type, contrasted with onscreen computer output	<code>machine_name%</code> su Password:
<i>AaBbCc123</i>	Command-line placeholder: replace with a real name or value	The command to remove a file is <code>rm filename</code> .
<i>AaBbCc123</i>	Book titles, new terms, and terms to be emphasized	Read Chapter 6 in the <i>User's Guide</i> . These are called <i>class</i> options. Do <i>not</i> save the file. (Emphasis sometimes appears in bold online.)

Shell Prompts in Command Examples

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

TABLE P-2 Shell Prompts

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

PART I Resource Management

This part introduces Solaris 10 Resource Manager, which enables you to control how applications use available system resources.

Introduction to Solaris 10 Resource Manager

Resource management functionality is a component of the Solaris™ Container environment. Resource management enables you to control how applications use available system resources. You can do the following:

- Allocate computing resources, such as processor time
- Monitor how the allocations are being used, then adjust the allocations as necessary
- Generate extended accounting information for analysis, billing, and capacity planning

This chapter covers the following topics.

- [“What’s New in Resource Management?”](#) on page 23
- [“Resource Management Overview”](#) on page 26
- [“When to Use Resource Management”](#) on page 29
- [“Setting Up Resource Management \(Task Map\)”](#) on page 30

What’s New in Resource Management?

This section describes new resource management features in the Solaris 10 release.

Improvements to Project Database and Resource Control Commands

Enhancements include the following:

- Scaled value and unit modifier support for resource control values and commands
- Improved validation and easier manipulation of the project attributes field

- Revised output format and new options for the `prctl` and `projects` commands
- Ability to set user's default project through the `useradd` command and modify information by using the `usermod` and `passmgmt` commands

For more information, see the following man pages:

- `passmgmt(1M)`
- `projadd(1M)`
- `projmod(1M)`
- `useradd(1M)`
- `usermod(1M)`

Also see [Chapter 2](#) and [Chapter 6](#).

New Resource Controls

The following set of resource controls replaces the System V interprocess communication (IPC) `/etc/system` tunables:

- `project.max-shm-ids`
- `project.max-msg-ids`
- `project.max-sem-ids`
- `project.max-shm-memory`
- `process.max-sem-nsems`
- `process.max-sem-ops`
- `process.max-msg-qbytes`

The following event port resource controls have been added:

- `project.max-device-locked-memory`
- `project.max-port-ids`
- `process.max-port-events`

The following cryptographic resource control has been added:

- `project.max-crypto-memory`

The following additional resource controls have been added:

- `project.max-lwps`
- `project.max-tasks`
- `project.max-contracts`

For more information, see [“Available Resource Controls”](#) on page 74.

Resource Pools Changes

Resource pools now provide a mechanism for dynamically adjusting each pool's resource allocation in response to system events and application load changes. Dynamic resource pools simplify and reduce the number of decisions required from an administrator. Adjustments are automatically made to preserve the system performance goals specified by an administrator. For more information, see [Chapter 12](#).

You can now use the `projmod` command to set the `project.pool` attribute in the `/etc/project` file.

New Extended Accounting Features

The following new features have been added:

- An Internet Protocol Quality of Service (IPQoS) flow accounting module, described in "Using Flow Accounting and Statistics Gathering (Tasks)" in *System Administration Guide: IP Services*, allows you to capture network flow information on a system.
- A Practical Extraction and Report Language (Perl) interface to `libexacct` is now available. The interface enables you to develop customized reporting and extraction scripts. See "[Perl Interface to libexacct](#)" on page 59.

Ability to Regulate Physical Memory Consumption by Processes

The resource capping daemon `rcapd` enables you to regulate physical memory consumption by processes running in projects that have resource caps defined. For more information, see [Chapter 10](#).

You can now use the `projmod` command to set the `rcap.max-rss` attribute in the `/etc/project` file.

Interaction With Solaris Zones

Resource management features can now be used with zones to further refine the application environment. Interactions between these features and zones are described in applicable sections in this guide.

Resource Management Overview

Modern computing environments have to provide a flexible response to the varying workloads that are generated by different applications on a system. A *workload* is an aggregation of all processes of an application or group of applications. If resource management features are not used, the Solaris Operating System responds to workload demands by adapting to new application requests dynamically. This default response generally means that all activity on the system is given equal access to resources. Solaris resource management features enable you to treat workloads individually. You can do the following:

- Restrict access to a specific resource
- Offer resources to workloads on a preferential basis
- Isolate workloads from each another

The ability to minimize cross-workload performance compromises, along with the facilities that monitor resource usage and utilization, is referred to as *resource management*. Resource management is implemented through a collection of algorithms. The algorithms handle the series of capability requests that an application presents in the course of its execution.

Resource management facilities permit you to modify the default behavior of the operating system with respect to different workloads. *Behavior* primarily refers to the set of decisions that are made by operating system algorithms when an application presents one or more resource requests to the system. You can use resource management facilities to do the following:

- Deny resources or prefer one application over another for a larger set of allocations than otherwise permitted
- Treat certain allocations collectively instead of through isolated mechanisms

The implementation of a system configuration that uses the resource management facilities can serve several purposes. You can do the following:

- Prevent an application from consuming resources indiscriminately
- Change an application's priority based on external events
- Balance resource guarantees to a set of applications against the goal of maximizing system utilization

When planning a resource-managed configuration, key requirements include the following:

- Identifying the competing workloads on the system
- Distinguishing those workloads that are not in conflict from those workloads with performance requirements that compromise the primary workloads

After you identify cooperating and conflicting workloads, you can create a resource configuration that presents the least compromise to the service goals of the business, within the limitations of the system's capabilities.

Effective resource management is enabled in the Solaris system by offering control mechanisms, notification mechanisms, and monitoring mechanisms. Many of these capabilities are provided through enhancements to existing mechanisms such as the `proc(4)` file system, processor sets, and scheduling classes. Other capabilities are specific to resource management. These capabilities are described in subsequent chapters.

Resource Classifications

A resource is any aspect of the computing system that can be manipulated with the intent to change application behavior. Thus, a resource is a capability that an application implicitly or explicitly requests. If the capability is denied or constrained, the execution of a robustly written application proceeds more slowly.

Classification of resources, as opposed to identification of resources, can be made along a number of axes. The axes could be implicitly requested as opposed to explicitly requested, time-based, such as CPU time, compared to time-independent, such as assigned CPU shares, and so forth.

Generally, scheduler-based resource management is applied to resources that the application can implicitly request. For example, to continue execution, an application implicitly requests additional CPU time. To write data to a network socket, an application implicitly requests bandwidth. Constraints can be placed on the aggregate total use of an implicitly requested resource.

Additional interfaces can be presented so that bandwidth or CPU service levels can be explicitly negotiated. Resources that are explicitly requested, such as a request for an additional thread, can be managed by constraint.

Resource Management Control Mechanisms

The three types of control mechanisms that are available in the Solaris Operating System are constraints, scheduling, and partitioning.

Constraint Mechanisms

Constraints allow the administrator or application developer to set bounds on the consumption of specific resources for a workload. With known bounds, modeling resource consumption scenarios becomes a simpler process. Bounds can also be used to control ill-behaved applications that would otherwise compromise system performance or availability through unregulated resource requests.

Constraints do present complications for the application. The relationship between the application and the system can be modified to the point that the application is no longer able to function. One approach that can mitigate this risk is to gradually narrow the constraints on applications with unknown resource behavior. The resource controls feature discussed in [Chapter 6](#) provides a constraint mechanism. Newer applications can be written to be aware of their resource constraints, but not all application writers will choose to do this.

Scheduling Mechanisms

Scheduling refers to making a sequence of allocation decisions at specific intervals. The decision that is made is based on a predictable algorithm. An application that does not need its current allocation leaves the resource available for another application's use. Scheduling-based resource management enables full utilization of an undercommitted configuration, while providing controlled allocations in a critically committed or overcommitted scenario. The underlying algorithm defines how the term "controlled" is interpreted. In some instances, the scheduling algorithm might guarantee that all applications have some access to the resource. The fair share scheduler (FSS) described in [Chapter 8](#) manages application access to CPU resources in a controlled way.

Partitioning Mechanisms

Partitioning is used to bind a workload to a subset of the system's available resources. This binding guarantees that a known amount of resources is always available to the workload. The resource pools functionality that is described in [Chapter 12](#) enables you to limit workloads to specific subsets of the machine.

Configurations that use partitioning can avoid system-wide overcommitment. However, in avoiding this overcommitment, the ability to achieve high utilizations can be reduced. A reserved group of resources, such as processors, is not available for use by another workload when the workload bound to them is idle.

Resource Management Configuration

Portions of the resource management configuration can be placed in a network name service. This feature allows the administrator to apply resource management constraints across a collection of machines, rather than on an exclusively per-machine basis. Related work can share a common identifier, and the aggregate usage of that work can be tabulated from accounting data.

Resource management configuration and workload-oriented identifiers are described more fully in [Chapter 2](#). The extended accounting facility that links these identifiers with application resource usage is described in [Chapter 4](#).

When to Use Resource Management

Use resource management to ensure that your applications have the required response times.

Resource management can also increase resource utilization. By categorizing and prioritizing usage, you can effectively use reserve capacity during off-peak periods, often eliminating the need for additional processing power. You can also ensure that resources are not wasted because of load variability.

Server Consolidation

Resource management is ideal for environments that consolidate a number of applications on a single server.

The cost and complexity of managing numerous machines encourages the consolidation of several applications on larger, more scalable servers. Instead of running each workload on a separate system, with full access to that system's resources, you can use resource management software to segregate workloads within the system. Resource management enables you to lower overall total cost of ownership by running and controlling several dissimilar applications on a single Solaris system.

If you are providing Internet and application services, you can use resource management to do the following:

- Host multiple web servers on a single machine. You can control the resource consumption for each web site and you can protect each site from the potential excesses of other sites.
- Prevent a faulty common gateway interface (CGI) script from exhausting CPU resources.
- Stop an incorrectly behaving application from leaking all available virtual memory.
- Ensure that one customer's applications are not affected by another customer's applications that run at the same site.
- Provide differentiated levels or classes of service on the same machine.
- Obtain accounting information for billing purposes.

Supporting a Large or Varied User Population

Use resource management features in any system that has a large, diverse user base, such as an educational institution. If you have a mix of workloads, the software can be configured to give priority to specific projects.

For example, in large brokerage firms, traders intermittently require fast access to execute a query or to perform a calculation. Other system users, however, have more consistent workloads. If you allocate a proportionately larger amount of processing power to the traders' projects, the traders have the responsiveness that they need.

Resource management is also ideal for supporting thin-client systems. These platforms provide stateless consoles with frame buffers and input devices, such as smart cards. The actual computation is done on a shared server, resulting in a timesharing type of environment. Use resource management features to isolate the users on the server. Then, a user who generates excess load does not monopolize hardware resources and significantly impact others who use the system.

Setting Up Resource Management (Task Map)

The following task map provides a high-level overview of the steps that are involved in setting up resource management on your system.

Task	Description	For Instructions
Identify the workloads on your system.	Review project entries in either the <code>/etc/project</code> file or in the NIS map or LDAP directory service.	"project Database" on page 35
Prioritize the workloads on your system.	Determine which applications are critical. These workloads might require preferential access to resources.	Refer to your business service goals.
Monitor real-time activity on your system.	Use performance tools to view the current resource consumption of workloads that are running on your system. You can then evaluate whether you must restrict access to a given resource or isolate particular workloads from other workloads.	"Monitoring by System" on page 186 and <code>cpustat(1M)</code> , <code>iostat(1M)</code> , <code>mpstat(1M)</code> , <code>prstat(1M)</code> , <code>sar(1)</code> , and <code>vmstat(1M)</code> man pages

Task	Description	For Instructions
Make temporary modifications to the workloads that are running on your system.	To determine which values can be altered, refer to the resource controls that are available in the Solaris system. You can update the values from the command line while the task or process is running.	“Available Resource Controls” on page 74, “Global and Local Actions on Resource Control Values” on page 79, “Temporarily Updating Resource Control Values on a Running System” on page 83 and <code>rct1adm(1M)</code> and <code>prctl(1)</code> man pages.
Set resource controls and project attributes for every project entry in the <code>project</code> database or naming service project database.	Each project entry in the <code>/etc/project</code> file or the naming service project database can contain one or more resource controls or attributes. Resource controls constrain tasks and processes attached to that project. For each threshold value that is placed on a resource control, you can associate one or more actions to be taken when that value is reached. You can set resource controls by using the command-line interface. Certain configuration parameters can also be set by using the Solaris Management Console.	“ <code>project</code> Database” on page 35, “Local <code>/etc/project</code> File Format” on page 36, “Available Resource Controls” on page 74, “Global and Local Actions on Resource Control Values” on page 79, and Chapter 8
Place an upper bound on the resource consumption of physical memory by collections of processes attached to a project.	The resource cap enforcement daemon will enforce the physical memory resource cap defined for the project’s <code>rcap.max-rss</code> attribute in the <code>/etc/project</code> file.	“ <code>project</code> Database” on page 35 and Chapter 10
Create resource pool configurations.	Resource pools provide a way to partition system resources, such as processors, and maintain those partitions across reboots. You can add one <code>project.pool</code> attribute to each entry in the <code>/etc/project</code> file.	“ <code>project</code> Database” on page 35 and Chapter 12
Make the fair share scheduler (FSS) your default system scheduler.	Ensure that all user processes in either a single CPU system or a processor set belong to the same scheduling class.	“Configuring the FSS” on page 109 and <code>dispadm(1M)</code> man page

Task	Description	For Instructions
<p>Activate the extended accounting facility to monitor and record resource consumption on a task or process basis.</p>	<p>Use extended accounting data to assess current resource controls and to plan capacity requirements for future workloads. Aggregate usage on a system-wide basis can be tracked. To obtain complete usage statistics for related workloads that span more than one system, the project name can be shared across several machines.</p>	<p>“How to Activate Extended Accounting for Processes, Tasks, and Flows” on page 64 and <code>acctadm(1M)</code> man page</p>
<p>(Optional) If you need to make additional adjustments to your configuration, you can continue to alter the values from the command line. You can alter the values while the task or process is running.</p>	<p>Modifications to existing tasks can be applied on a temporary basis without restarting the project. Tune the values until you are satisfied with the performance. Then, update the current values in the <code>/etc/project</code> file or in the naming service project database.</p>	<p>“Temporarily Updating Resource Control Values on a Running System” on page 83 and <code>rctladm(1M)</code> and <code>prctl(1)</code> man pages</p>
<p>(Optional) Capture extended accounting data.</p>	<p>Write extended accounting records for active processes and active tasks. The files that are produced can be used for planning, chargeback, and billing purposes. There is also a Practical Extraction and Report Language (Perl) interface to <code>libexacct</code> that enables you to develop customized reporting and extraction scripts.</p>	<p><code>wracct(1M)</code> man page and “Perl Interface to libexacct” on page 59</p>

Projects and Tasks (Overview)

This chapter discusses the *project* and *task* facilities of Solaris resource management. Projects and tasks are used to label workloads and separate them from one another.

The following topics are covered in this chapter:

- “Project and Task Facilities” on page 33
- “Project Identifiers” on page 34
- “Task Identifiers” on page 39
- “Commands Used With Projects and Tasks” on page 40

To use the projects and tasks facilities, see [Chapter 3](#).

Project and Task Facilities

To optimize workload response, you must first be able to identify the workloads that are running on the system you are analyzing. This information can be difficult to obtain by using either a purely process-oriented or a user-oriented method alone. In the Solaris system, you have two additional facilities that can be used to separate and identify workloads: the project and the task. The *project* provides a network-wide administrative identifier for related work. The *task* collects a group of processes into a manageable entity that represents a workload component.

The controls specified in the `project` name service database are set on the process, task, and project. Since process and task controls are inherited across `fork` and `settaskid` system calls, all processes and tasks that are created within the project inherit these controls. For information on these system calls, see the `fork(2)` and `settaskid(2)` man pages.

Based on their project or task membership, running processes can be manipulated with standard Solaris commands. The extended accounting facility can report on both process usage and task usage, and tag each record with the governing project

identifier. This process enables offline workload analysis to be correlated with online monitoring. The project identifier can be shared across multiple machines through the project name service database. Thus, the resource consumption of related workloads that run on (or span) multiple machines can ultimately be analyzed across all of the machines.

Project Identifiers

The project identifier is an administrative identifier that is used to identify related work. The project identifier can be thought of as a workload tag equivalent to the user and group identifiers. A user or group can belong to one or more projects. These projects can be used to represent the workloads in which the user (or group of users) is allowed to participate. This membership can then be the basis of chargeback that is based on, for example, usage or initial resource allocations. Although a user must be assigned to a default project, the processes that the user launches can be associated with any of the projects of which that user is a member.

Determining a User's Default Project

To log in to the system, a user must be assigned a default project. A user is automatically a member of that default project, even if the user is not in the user or group list specified in that project.

Because each process on the system possesses project membership, an algorithm to assign a default project to the login or other initial process is necessary. The algorithm is documented in the man page `getproject(3C)`. The system follows ordered steps to determine the default project. If no default project is found, the user's login, or request to start a process, is denied.

The system sequentially follows these steps to determine a user's default project:

1. If the user has an entry with a `project` attribute defined in the `/etc/user_attr` extended user attributes database, then the value of the `project` attribute is the default project. See the `user_attr(4)` man page.
2. If a project with the name `user.user-id` is present in the `project` database, then that project is the default project. See the `project(4)` man page for more information.
3. If a project with the name `group.group-name` is present in the `project` database, where `group-name` is the name of the default group for the user, as specified in the `passwd` file, then that project is the default project. For information on the `passwd` file, see the `passwd(4)` man page.
4. If the special project `default` is present in the `project` database, then that project is the default project.

This logic is provided by the `getdefaultproj()` library function. See the `getproject(3PROJECT)` man page for more information.

Setting User Attributes With the `useradd`, `usermod`, and `passmgmt` Commands

You can use the following commands with the `-K` option and a *key=value* pair to set user attributes in local files :

<code>passmgmt</code>	modify user information
<code>useradd</code>	set default project for user
<code>usermod</code>	modify user information

Local files can include the following:

- `/etc/group`
- `/etc/passwd`
- `/etc/project`
- `/etc/shadow`
- `/etc/user_attr`

If a network naming service such as NIS is being used to supplement the local file with additional entries, these commands cannot change information supplied by the network name service. However, the commands do verify the following against the external *naming service database*:

- Uniqueness of the user name (or role)
- Uniqueness of the user ID
- Existence of any group names specified

For more information, see the `passmgmt(1M)`, `useradd(1M)`, `usermod(1M)`, and `user_attr(4)` man pages.

project Database

You can store project data in a local file, in a Network Information Service (NIS) project map, or in a Lightweight Directory Access Protocol (LDAP) directory service. The `/etc/project` file or naming service is used at login and by all requests for account management by the pluggable authentication module (PAM) to bind a user to a default project.

Note – Updates to entries in the project database, whether to the `/etc/project` file or to a representation of the database in a network naming service, are not applied to currently active projects. The updates are applied to new tasks that join the project when either the `login` or the `newtask` command is used. For more information, see the `login(1)` and `newtask(1)` man pages.

PAM Subsystem

Operations that change or set identity include logging in to the system, invoking an `rcp` or `rsh` command, using `ftp`, or using `su`. When an operation involves changing or setting an identity, a set of configurable modules is used to provide authentication, account management, credentials management, and session management.

The account management PAM module for projects is documented in the `pam_projects(5)` man page. For an overview of PAM, see “Using PAM” in *System Administration Guide: Security Services*.

Naming Services Configuration

Resource management supports naming service project databases. The location where the project database is stored is defined in the `/etc/nsswitch.conf` file. By default, `files` is listed first, but the sources can be listed in any order.

```
project: files [nis] [ldap]
```

If more than one source for project information is listed, the `nsswitch.conf` file directs the routine to start searching for the information in the first source listed, and then search subsequent sources.

For more information about the `/etc/nsswitch.conf` file, see “The Name Service Switch” in *System Administration Guide: Naming and Directory Services (DNS, NIS, and LDAP)* and `nsswitch.conf(4)`.

Local `/etc/project` File Format

If you select `files` as your project database source in the `nsswitch.conf` file, the `login` process searches the `/etc/project` file for project information. See the `projects(1)` and `project(4)` man pages for more information.

The `project` file contains a one-line entry of the following form for each project recognized by the system:

projname:projid:comment:user-list:group-list:attributes

The fields are defined as follows:

- projname* The name of the project. The name must be a string that consists of alphanumeric characters, underline (`_`) characters, hyphens (`-`), and periods (`.`). The period, which is reserved for projects with special meaning to the operating system, can only be used in the names of default projects for users. *projname* cannot contain colons (`:`) or newline characters.
- projid* The project's unique numerical ID (PROJID) within the system. The maximum value of the *projid* field is `UID_MAX` (2147483647).
- comment* A description of the project.
- user-list* A comma-separated list of users who are allowed in the project.
- Wildcards can be used in this field. An asterisk (`*`) allows all users to join the project. An exclamation point followed by an asterisk (`!*`) excludes all users from the project. An exclamation mark (`!`) followed by a user name excludes the specified user from the project.
- group-list* A comma-separated list of groups of users who are allowed in the project.
- Wildcards can be used in this field. An asterisk (`*`) allows all groups to join the project. An exclamation point followed by an asterisk (`!*`) excludes all groups from the project. An exclamation mark (`!`) followed by a group name excludes the specified group from the project.
- attributes* A semicolon-separated list of name-value pairs, such as resource controls (see [Chapter 6](#)). *name* is an arbitrary string that specifies the object-related attribute, and *value* is the optional value for that attribute.
- `name [=value]`
- In the name-value pair, names are restricted to letters, digits, underscores, and periods. A period is conventionally used as a separator between the categories and subcategories of the resource control (rctl). The first character of an attribute name must be a letter. The name is case sensitive.
- Values can be structured by using commas and parentheses to establish precedence.
- A semicolon is used to separate name-value pairs. A semicolon cannot be used in a value definition. A colon is used to separate project fields. A colon cannot be used in a value definition.

Note – Routines that read this file halt if they encounter a malformed entry. Any projects that are specified after the incorrect entry are not assigned.

This example shows the default `/etc/project` file:

```
system:0:System:::
user.root:1:Super-User:::
noproject:2:No Project:::
default:3:::
group.staff:10:::
```

This example shows the default `/etc/project` file with project entries added at the end:

```
system:0:System:::
user.root:1:Super-User:::
noproject:2:No Project:::
default:3:::
group.staff:10:::
user.ml:2424:Lyle Personal:::
booksite:4113:Book Auction Project:ml,mp,jtd,kjh:::
```

You can also add resource controls and attributes to the `/etc/project` file:

- To add resource controls for a project, see [“Setting Resource Controls”](#) on page 86.
- To define a physical memory resource cap for a project using the resource capping daemon described in `rcapd(1M)`, see [“Attribute to Limit Physical Memory Usage”](#) on page 115.
- To add a project `.pool` attribute to a project’s entry, see [“Creating the Configuration”](#) on page 176.

Project Configuration for NIS

If you are using NIS, you can specify in the `/etc/nsswitch.conf` file to search the NIS project maps for projects:

```
project: nis files
```

The NIS maps, either `project.byname` or `project.bynumber`, have the same form as the `/etc/project` file:

```
projname:projid:comment:user-list:group-list:attributes
```

For more information, see [“Network Information Service \(NIS\) \(Overview\)”](#) in *System Administration Guide: Naming and Directory Services (DNS, NIS, and LDAP)*.

Project Configuration for LDAP

If you are using LDAP, you can specify in the `/etc/nsswitch.conf` file to search the LDAP `project` database for projects:

```
project: ldap files
```

For more information about LDAP, see “Introduction to LDAP Naming Services (Overview/Reference)” in *System Administration Guide: Naming and Directory Services (DNS, NIS, and LDAP)*. For more information about the schema for project entries in an LDAP database, see “Solaris Schemas” in *System Administration Guide: Naming and Directory Services (DNS, NIS, and LDAP)*.

Task Identifiers

Each successful login into a project creates a new *task* that contains the login process. The task is a process collective that represents a set of work over time. A task can also be viewed as a *workload component*. Each task is automatically assigned a task ID.

Each process is a member of one task, and each task is associated with one project.

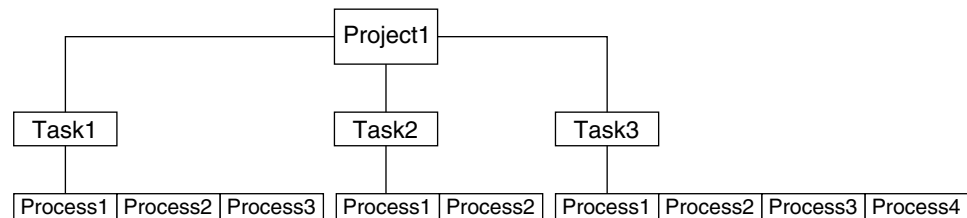


FIGURE 2-1 Project and Task Tree

All operations on process groups, such as signal delivery, are also supported on tasks. You can also bind a task to a *processor set* and set a scheduling priority and class for a task, which modifies all current and subsequent processes in the task.

A task is created whenever a project is joined. The following actions, commands, and functions create tasks:

- `login`
- `cron`
- `newtask`
- `setproject`
- `su`

You can create a finalized task by using one of the following methods. All further attempts to create new tasks will fail.

- You can use the `newtask` command with the `-F` option.
- You can set the `task.final` attribute on a project in the project naming service database. All tasks created in that project by `setproject` have the `TASK_FINAL` flag.

For more information, see the `login(1)`, `newtask(1)`, `cron(1M)`, `su(1M)`, and `setproject(3PROJECT)` man pages.

The extended accounting facility can provide accounting data for processes. The data is aggregated at the task level.

Commands Used With Projects and Tasks

The commands that are shown in the following table provide the primary administrative interface to the project and task facilities.

Man Page Reference	Description
<code>projects(1)</code>	Displays project memberships for users. Lists projects from project database. Prints information on given projects. If no project names are supplied, information is displayed for all projects. Use the <code>projects</code> command with the <code>-l</code> option to print verbose output.
<code>newtask(1)</code>	Executes the user's default shell or specified command, placing the execution command in a new task that is owned by the specified project. <code>newtask</code> can also be used to change the task and the project binding for a running process. Use with the <code>-F</code> option to create a finalized task.
<code>passmgmt(1M)</code>	Updates information in the password files. Use with the <code>-K key=value</code> option to add to user attributes or replace user attributes in local files.

Man Page Reference	Description
projadd(1M)	<p>Adds a new project entry to the <code>/etc/project</code> file. The <code>projadd</code> command creates a project entry only on the local system. <code>projadd</code> cannot change information that is supplied by the network naming service.</p> <p>Can be used to edit project files other than the default file, <code>/etc/project</code>. Provides syntax checking for <code>project</code> file. Validates and edits project attributes. Supports scaled values.</p>
projmod(1M)	<p>Modifies information for a project on the local system. <code>projmod</code> cannot change information that is supplied by the network naming service. However, the command does verify the uniqueness of the project name and project ID against the external naming service.</p> <p>Can be used to edit project files other than the default file, <code>/etc/project</code>. Provides syntax checking for <code>project</code> file. Validates and edits project attributes. Can be used to add a new attribute, add values to an attribute, or remove an attribute. Supports scaled values.</p>
projdel(1M)	<p>Deletes a project from the local system. <code>projdel</code> cannot change information that is supplied by the network naming service.</p>
useradd(1M)	<p>Adds default project definitions to the local files. Use with the <code>-K key=value</code> option to add or replace user attributes.</p>
userdel(1M)	<p>Deletes a user's account from the local file.</p>
usermod(1M)	<p>Modifies a user's login information on the system. Use with the <code>-K key=value</code> option to add or replace user attributes.</p>

Administering Projects and Tasks

This chapter describes how to use the project and task facilities of Solaris resource management.

The following topics are covered.

- “Sample Commands and Command Options” on page 44
- “Administering Projects” on page 47

For an overview of the projects and tasks facilities, see [Chapter 2](#).

Note – If you are using these facilities in a zones environment, only processes in the same zone will be visible through system call interfaces that take process IDs when these commands are run in a non-global zone.

Administering Projects and Tasks (Task Map)

Task	Description	For Instructions
View examples of commands and options used with projects and tasks.	Display task and project IDs, display various statistics for processes and projects that are currently running on your system.	“Sample Commands and Command Options” on page 44

Task	Description	For Instructions
Define a project.	Add a project entry to the <code>/etc/project</code> file and alter values for that entry.	“How to Define a Project and View the Current Project” on page 47
Delete a project.	Remove a project entry from the <code>/etc/project</code> file.	“How to Delete a Project From the <code>/etc/project</code> File” on page 49
Validate the <code>project</code> file or project database.	Check the syntax of the <code>/etc/project</code> file or verify the uniqueness of the project name and project ID against the external naming service.	“How to Validate the Contents of the <code>/etc/project</code> File” on page 50
Obtain project membership information.	Display the current project membership of the invoking process.	“How to Obtain Project Membership Information” on page 51
Create a new task.	Create a new task in a particular project by using the <code>newtask</code> command.	“How to Create a New Task” on page 51
Associate a running process with a different task and project.	Associate a process number with a new task ID in a specified project.	“How to Move a Running Process Into a New Task” on page 51
Add and work with project attributes.	Use the project database administration commands to add, edit, validate, and remove project attributes.	“Editing and Validating Project Attributes” on page 52

Sample Commands and Command Options

This section provides examples of commands and options used with projects and tasks.

Command Options Used With Projects and Tasks

`ps` Command

Use the `ps` command with the `-o` option to display task and project IDs. For example, to view the project ID, type the following:

```
# ps -o user,pid,uid,projid
USER PID  UID  PROJID
jtd  89430 124  4113
```

id Command

Use the `id` command with the `-p` option to print the current project ID in addition to the user and group IDs. If the `user` operand is provided, the project associated with that user's normal login is printed:

```
# id -p
uid=124(jtd) gid=10(staff) projid=4113(booksite)
```

pgrep and pkill Commands

To match only processes with a project ID in a specific list, use the `pgrep` and `pkill` commands with the `-J` option:

```
# pgrep -J projidlist
# pkill -J projidlist
```

To match only processes with a task ID in a specific list, use the `pgrep` and `pkill` commands with the `-T` option:

```
# pgrep -T taskidlist
# pkill -T taskidlist
```

prstat Command

To display various statistics for processes and projects that are currently running on your system, use the `prstat` command with the `-J` option:

```
% prstat -J
      PID USERNAME  SIZE  RSS STATE  PRI NICE   TIME  CPU PROCESS/NLWP
21634 jtd      5512K 4848K cpu0   44  0   0:00.00 0.3% prstat/1
   324 root        29M   75M sleep   59  0   0:08.27 0.2% Xsun/1
15497 jtd       48M   41M sleep   49  0   0:08.26 0.1% adeptedit/1
   328 root     2856K 2600K sleep   58  0   0:00.00 0.0% mibiisa/11
  1979 jtd     1568K 1352K sleep   49  0   0:00.00 0.0% csh/1
  1977 jtd     7256K 5512K sleep   49  0   0:00.00 0.0% dtterm/1
   192 root     3680K 2856K sleep   58  0   0:00.36 0.0% automountd/5
  1845 jtd       24M   22M sleep   49  0   0:00.29 0.0% dtmail/11
  1009 jtd     9864K 8384K sleep   49  0   0:00.59 0.0% dtwm/8
   114 root     1640K  704K sleep   58  0   0:01.16 0.0% in.routed/1
   180 daemon  2704K 1944K sleep   58  0   0:00.00 0.0% statd/4
   145 root     2120K 1520K sleep   58  0   0:00.00 0.0% ypbind/1
   181 root     1864K 1336K sleep   51  0   0:00.00 0.0% lockd/1
   173 root     2584K 2136K sleep   58  0   0:00.00 0.0% inetd/1
   135 root     2960K 1424K sleep    0  0   0:00.00 0.0% keyserv/4
```

PROJID	NPROC	SIZE	RSS	MEMORY	TIME	CPU	PROJECT
10	52	400M	271M	68%	0:11.45	0.4%	booksite
0	35	113M	129M	32%	0:10.46	0.2%	system

Total: 87 processes, 205 lwps, load averages: 0.05, 0.02, 0.02

To display various statistics for processes and tasks that are currently running on your system, use the `prstat` command with the `-T` option:

```
% prstat -T
  PID USERNAME  SIZE  RSS STATE PRI NICE      TIME  CPU PROCESS/NLWP
23023 root         26M   20M sleep  59   0   0:03:18  0.6% Xsun/1
23476 jtd         51M   45M sleep  49   0   0:04:31  0.5% adeptedit/1
23432 jtd        6928K 5064K sleep  59   0   0:00:00  0.1% dtterm/1
28959 jtd         26M   18M sleep  49   0   0:00:18  0.0% .netscape.bin/1
23116 jtd        9232K 8104K sleep  59   0   0:00:27  0.0% dtwm/5
29010 jtd        5144K 4664K cpu0   59   0   0:00:00  0.0% prstat/1
  200 root        3096K 1024K sleep  59   0   0:00:00  0.0% lpsched/1
  161 root        2120K 1600K sleep  59   0   0:00:00  0.0% lockd/2
  170 root        5888K 4248K sleep  59   0   0:03:10  0.0% automountd/3
  132 root        2120K 1408K sleep  59   0   0:00:00  0.0% ypbind/1
  162 daemon      2504K 1936K sleep  59   0   0:00:00  0.0% statd/2
  146 root        2560K 2008K sleep  59   0   0:00:00  0.0% inetd/1
  122 root        2336K 1264K sleep  59   0   0:00:00  0.0% keyserv/2
  119 root        2336K 1496K sleep  59   0   0:00:02  0.0% rpcbind/1
  104 root        1664K  672K sleep  59   0   0:00:03  0.0% in.rdisc/1
TASKID  NPROC  SIZE  RSS MEMORY  TIME  CPU PROJECT
  222    30  229M  161M   44%   0:05:54  0.6% group.staff
  223    1   26M   20M    5.3%   0:03:18  0.6% group.staff
   12    1   61M   33M    8.9%   0:00:31  0.0% group.staff
    1    33   85M   53M   14%   0:03:33  0.0% system
```

Total: 65 processes, 154 lwps, load averages: 0.04, 0.05, 0.06

Note – The `-J` and `-T` options cannot be used together.

Using cron and su With Projects and Tasks

cron Command

The `cron` command issues a `settaskid` to ensure that each `cron`, `at`, and `batch` job executes in a separate task, with the appropriate default project for the submitting user. The `at` and `batch` commands also capture the current project ID, which ensures that the project ID is restored when running an `at` job.

su Command

The `su` command joins the target user's default project by creating a new task, as part of simulating a login.

To switch the user's default project by using the `su` command, type the following:

```
# su user
```

Administering Projects

▼ How to Define a Project and View the Current Project

This example shows how to use the `projadd` command to add a project entry and the `projmod` command to alter that entry.

1. Become superuser or assume an equivalent role.

Roles contain authorizations and privileged commands. For more information about roles, see “Using the Solaris Management Tools With RBAC (Task Map)” in *System Administration Guide: Basic Administration*.

2. View the default `/etc/project` file on your system by using `projects -l`.

```
# projects -l
system:0:::
user.root:1:::
noproject:2:::
default:3:::
group.staff:10:::system
    projid : 0
    comment: ""
    users  : (none)
    groups : (none)
    attribs:
user.root
    projid : 1
    comment: ""
    users  : (none)
    groups : (none)
    attribs:
noproject
    projid : 2
    comment: ""
    users  : (none)
    groups : (none)
    attribs:
default
    projid : 3
    comment: ""
```

```
        users : (none)
        groups : (none)
        attribs:
group.staff
    projid : 10
    comment: ""
    users : (none)
    groups : (none)
    attribs:
```

3. Add a project with the name *booksite*. Assign the project to a user who is named *mark* with project ID number *4113*.

```
# projadd -U mark -p 4113 booksite
```

4. View the `/etc/project` file again.

```
# projects -l
system
    projid : 0
    comment: ""
    users : (none)
    groups : (none)
    attribs:
user.root
    projid : 1
    comment: ""
    users : (none)
    groups : (none)
    attribs:
noproject
    projid : 2
    comment: ""
    users : (none)
    groups : (none)
    attribs:
default
    projid : 3
    comment: ""
    users : (none)
    groups : (none)
    attribs:
group.staff
    projid : 10
    comment: ""
    users : (none)
    groups : (none)
    attribs:
booksite
    projid : 4113
    comment: ""
    users : mark
    groups : (none)
    attribs:
```


5. Add a comment that describes the project in the comment field.

```
# projmod -c 'Book Auction Project' booksite
```

6. View the changes in the `/etc/project` file.

```
# projects -l
system
    projid : 0
    comment: ""
    users  : (none)
    groups : (none)
    attribs:
user.root
    projid : 1
    comment: ""
    users  : (none)
    groups : (none)
    attribs:
noproject
    projid : 2
    comment: ""
    users  : (none)
    groups : (none)
    attribs:
default
    projid : 3
    comment: ""
    users  : (none)
    groups : (none)
    attribs:
group.staff
    projid : 10
    comment: ""
    users  : (none)
    groups : (none)
    attribs:
booksite
    projid : 4113
    comment: "Book Auction Project"
    users  : mark
    groups : (none)
    attribs:
```

To bind projects, tasks, and processes to a pool, see [“Binding to a Pool”](#) on page 171.

▼ How to Delete a Project From the `/etc/project` File

This example shows how to use the `projdel` command to delete a project.

1. Become superuser or assume an equivalent role.

Roles contain authorizations and privileged commands. For more information about roles, see “Using the Solaris Management Tools With RBAC (Task Map)” in *System Administration Guide: Basic Administration*.

2. Remove the project *booksite* by using the `projdel` command.

```
# projdel booksite
```

3. Display the `/etc/project` file.

```
# projects -l
system
    projid : 0
    comment: ""
    users  : (none)
    groups : (none)
    attribs:
user.root
    projid : 1
    comment: ""
    users  : (none)
    groups : (none)
    attribs:
noproject
    projid : 2
    comment: ""
    users  : (none)
    groups : (none)
    attribs:
default
    projid : 3
    comment: ""
    users  : (none)
    groups : (none)
    attribs:
group.staff
    projid : 10
    comment: ""
    users  : (none)
    groups : (none)
    attribs:
```

4. Log in as user *mark* and type `projects` to view the projects that are assigned to this user.

```
# su - mark
# projects
default
```

How to Validate the Contents of the `/etc/project` File

If no editing options are given, the `projmod` command validates the contents of the `project` file.

To validate a NIS map, type the following:

```
$ ypcat project | projmod -f -
```

To check the syntax of the `/etc/project` file, type the following:

```
$ projmod -s
```

How to Obtain Project Membership Information

Use the `id` command with the `-p` flag to display the current project membership of the invoking process.

```
$ id -p
uid=100(mark) gid=1(other) projid=3(default)
```

▼ How to Create a New Task

1. **Log in as a member of the destination project, *booksite*.**
2. **Create a new task in the *booksite* project by using the `newtask` command with the `-v` (verbose) option to obtain the system task ID.**

```
machine% newtask -v -p booksite
16
```

The execution of `newtask` creates a new task in the specified project, and places the user's default shell in this task.

3. **View the current project membership of the invoking process.**

```
# id -p
uid=100(mark) gid=1(other) projid=4113(booksite)
```

The process is now a member of the new project.

▼ How to Move a Running Process Into a New Task

This example shows how to associate a running process with a different task and new project. To perform this action, you must either be superuser, or be the owner of the process and be a member of the new project.

1. **Become superuser or assume an equivalent role.**

Roles contain authorizations and privileged commands. For more information about roles, see "Using the Solaris Management Tools With RBAC (Task Map)" in *System Administration Guide: Basic Administration*.

Note – If you are the owner of the process or a member of the new project, you can skip this step.

2. Obtain the process ID of the *book_catalog* process.

```
# pgrep book_catalog
8100
```

3. Associate process *8100* with a new task ID in the *booksite* project.

```
# newtask -v -p booksite -c 8100
17
```

The `-c` option specifies that `newtask` operate on the existing named process.

4. Confirm the task to process ID mapping.

```
# pgrep -T 17
8100
```

Editing and Validating Project Attributes

You can use the `projadd` and `projmod` project database administration commands to edit project attributes.

The `-K` option specifies a replacement list of attributes. Attributes are delimited by semicolons (;). If the `-K` option is used with the `-a` option, the attribute or attribute value is added. If the `-K` option is used with the `-r` option, the attribute or attribute value is replaced. If the `-K` option is used with the `-s` option, the attribute or attribute value is substituted.

How to Add Attributes and Attribute Values to Projects

Use the `projmod` command with the `-a` and `-K` options to add values to an attribute in the project *myproject*. If the attribute does not exist, it is created. The following line adds a `task.max-lwps` resource control attribute with no values. A task entering the project has only the system value for the attribute.

```
# projmod -a -K task.max-lwps myproject
```

You can then add a value to `task.max-lwps` in the project *myproject*. The value consists of a privilege level, a threshold value, and an action associated with reaching the threshold.

```
# projmod -a -K task.max-lwps=(priv,100,deny) myproject
```

Because resource controls can have multiple values, the new value is added to the existing list of values if the options are used again. The multiple values are separated by commas. For example, the following command:

```
# projmod -a -K task.max-lwps=(priv,1000,signal=KILL) myproject
```

Produces the `task.max-lwps` entry:

```
task.max-lwps=(priv,100,deny) ,(priv,1000,signal=KILL)
```

How to Remove Attributes and Attribute Values From Projects

To remove an attribute value from the resource control `task.max-lwps` in the project *myproject*, use the `projmod` command with the `-r` and `-K` options.

```
# projmod -r -K task.max-lwps=(priv,100,deny) myproject
```

If `task.max-lwps` has multiple values, such as:

```
task.max-lwps=(priv,100,deny) ,(priv,1000,signal=KILL)
```

The first matching value would be removed. The result would then be:

```
task.max-lwps=(priv,1000,signal=KILL)
```

The following command removes `task.max-lwps` and all of its values:

```
# projmod -r -K task.max-lwps myproject
```

How to Substitute Attributes and Attribute Values for Projects

To substitute a different value for the attribute `task.max-lwps` in the project *myproject*, use the `projmod` command with the `-s` and `-K` options. If the attribute does not exist, it is created.

The following command replaces the current `task.max-lwps` values with new values:

```
# projmod -s -K task.max-lwps=(priv,100,none) ,(priv,120,deny) myproject
```

The following command removes the current `task.max-lwps` values:

```
# projmod -s -K task.max-lwps myproject
```

Extended Accounting (Overview)

By using the project and task facilities that are described in [Chapter 2](#) to label and separate workloads, you can monitor resource consumption by each workload. You can use the *extended accounting* subsystem to capture a detailed set of resource consumption statistics on both processes and tasks.

The following topics are covered in this chapter.

- “Introduction to Extended Accounting” on page 55
- “How Extended Accounting Works” on page 56
- “Extended Accounting Configuration” on page 58
- “Commands Used With Extended Accounting” on page 59
- “Perl Interface to `libexacct`” on page 59

To begin using extended accounting, see “How to Activate Extended Accounting for Processes, Tasks, and Flows” on page 64.

Introduction to Extended Accounting

The extended accounting subsystem labels usage records with the project for which the work was done. You can also use extended accounting, in conjunction with the Internet Protocol Quality of Service (IPQoS) flow accounting module described in “Using Flow Accounting and Statistics Gathering (Tasks)” in *System Administration Guide: IP Services*, to capture network flow information on a system.

Before you can apply resource management mechanisms, you must first be able to characterize the resource consumption demands that various workloads place on a system. The extended accounting facility in the Solaris Operating System provides a flexible way to record system and network resource consumption on a task or process basis, or on the basis of selectors provided by the IPQoS `flowacct` module. For more information, see `ipqos(7IPP)`.

Unlike online monitoring tools, which enable you to measure system usage in real time, extended accounting enables you to examine historical usage. You can then make assessments of capacity requirements for future workloads.

With extended accounting data available, you can develop or purchase software for resource chargeback, workload monitoring, or capacity planning.

How Extended Accounting Works

The extended accounting facility in the Solaris Operating System uses a versioned, extensible file format to contain accounting data. Files that use this data format can be accessed or be created by using the API provided in the included library, `libexacct` (see `libexacct(3LIB)`). These files can then be analyzed on any platform with extended accounting enabled, and their data can be used for capacity planning and chargeback.

If extended accounting is active, statistics are gathered that can be examined by the `libexacct` API. `libexacct` allows examination of the `exacct` files either forward or backward. The API supports third-party files that are generated by `libexacct` as well as those files that are created by the kernel. There is a Practical Extraction and Report Language (Perl) interface to `libexacct` that enables you to develop customized reporting and extraction scripts. See [“Perl Interface to `libexacct`”](#) on page 59.

With extended accounting enabled, the task tracks the aggregate resource usage of its member processes. A task accounting record is written at task completion. Interim records on running processes and tasks can also be written. For more information on tasks, see [Chapter 2](#).

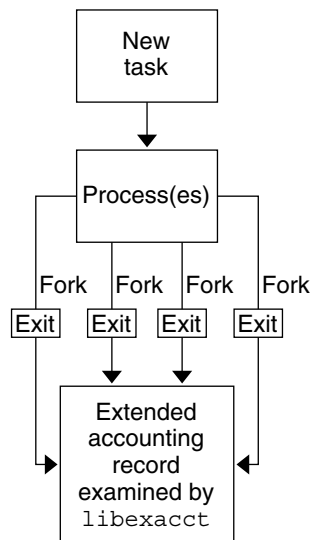


FIGURE 4-1 Task Tracking With Extended Accounting Activated

Extensible Format

The extended accounting format is substantially more extensible than the SunOS™ legacy system accounting software format (see “What is System Accounting?” in *System Administration Guide: Advanced Administration*). Extended accounting permits accounting metrics to be added and removed from the system between releases, and even during system operation.

Note – Both extended accounting and legacy system accounting software can be active on your system at the same time.

exacct Records and Format

Routines that allow `exacct` records to be created serve two purposes.

- To enable third-party `exacct` files to be created.
- To enable the creation of tagging records to be embedded in the kernel accounting file by using the `putacct` system call (see `getacct(2)`).

Note – The `putacct` system call is also available from the Perl interface.

The format permits different forms of accounting records to be captured without requiring that every change be an explicit version change. Well-written applications that consume accounting data must ignore records they do not understand.

The `libexacct` library converts and produces files in the `exacct` format. This library is the *only* supported interface to `exacct` format files.

Note – The `getacct`, `putacct`, and `wracct` system calls do not apply to flows. The kernel creates flow records and writes them to the file when IPQoS flow accounting is configured.

Using Extended Accounting in a Zones Environment

The extended accounting subsystem collects and reports information for the entire system (including non-global zones) when run in the global zone. The global administrator can also determine resource consumption on a per-zone basis. See “[Extended Accounting in a Zones Environment](#)” on page 307 for more information.

Extended Accounting Configuration

The `/etc/acctadm.conf` file contains the current extended accounting configuration. The file is edited through the `acctadm` interface, not by the user.

The directory `/var/adm/exacct` is the standard location for placing extended accounting data. You can use the `acctadm` command to specify a different location for the process and task accounting-data files. See `acctadm(1M)` for more information.

Commands Used With Extended Accounting

Command Reference	Description
<code>acctadm(1M)</code>	Modifies various attributes of the extended accounting facility, stops and starts extended accounting, and is used to select accounting attributes to track for processes, tasks, and flows.
<code>wracct(1M)</code>	Writes extended accounting records for active processes and active tasks.
<code>lastcomm(1)</code>	Displays previously invoked commands. <code>lastcomm</code> can consume either standard accounting-process data or extended-accounting process data.

For information on commands that are associated with tasks and projects, see [“Sample Commands and Command Options”](#) on page 44. For information on IPQoS flow accounting, see `ipqosconf(1M)`.

Perl Interface to `libexacct`

The Perl interface allows you to create Perl scripts that can read the accounting files produced by the `exacct` framework. You can also create Perl scripts that write `exacct` files.

The interface is functionally equivalent to the underlying C API. When possible, the data obtained from the underlying C API is presented as Perl data types. This feature makes accessing the data easier and it removes the need for `buffer pack` and `unpack` operations. Moreover, all memory management is performed by the Perl library.

The various project, task, and `exacct`-related functions are separated into groups. Each group of functions is located in a separate Perl module. Each module begins with the Sun standard `Sun::Solaris::` Perl package prefix. All of the classes provided by the Perl `exacct` library are found under the `Sun::Solaris::Exacct` module.

The underlying `libexacct(3LIB)` library provides operations on `exacct` format files, catalog tags, and `exacct` objects. `exacct` objects are subdivided into two types:

- Items, which are single-data values (scalars)
- Groups, which are lists of Items

The following table summarizes each of the modules.

Module (should not contain spaces)	Description	For More Information
Sun::Solaris::Project	This module provides functions to access the project manipulation functions <code>getprojid(2)</code> , <code>endprojent(3PROJECT)</code> , <code>fgetprojent(3PROJECT)</code> , <code>getdefaultproj(3PROJECT)</code> , <code>getprojbyid(3PROJECT)</code> , <code>getprojbyname(3PROJECT)</code> , <code>getprojent(3PROJECT)</code> , <code>getprojidbyname(3PROJECT)</code> , <code>inproj(3PROJECT)</code> , <code>project_walk(3PROJECT)</code> , <code>setproject(3PROJECT)</code> , and <code>setprojent(3PROJECT)</code> .	Project(3PERL)
Sun::Solaris::Task	This module provides functions to access the task manipulation functions <code>gettaskid(2)</code> and <code>settaskid(2)</code> .	Task(3PERL)
Sun::Solaris::Exacct	This module is the top-level <code>exacct</code> module. This module provides functions to access the <code>exacct</code> -related system calls <code>getacct(2)</code> , <code>putacct(2)</code> , and <code>wracct(2)</code> . This module also provides functions to access the <code>libexacct(3LIB)</code> library function <code>ea_error(3EXACCT)</code> . Constants for all of the <code>exacct</code> <code>EO_*</code> , <code>EW_*</code> , <code>EXR_*</code> , <code>P_*</code> , and <code>TASK_*</code> macros are also provided in this module.	Exacct(3PERL)
Sun::Solaris::Exacct::Catalog	This module provides object-oriented methods to access the bitfields in an <code>exacct</code> catalog tag. This module also provides access to the constants for the <code>EXC_*</code> , <code>EXD_*</code> , and <code>EXD_*</code> macros.	Exacct::Catalog(3PERL)
Sun::Solaris::Exacct::File	This module provides object-oriented methods to access the <code>libexacct</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> , and <code>ea_write_object(3EXACCT)</code> .	Exacct::File(3PERL)

Module (should not contain spaces)	Description	For More Information
Sun::Solaris::Exacct::Object	This module provides object-oriented methods to access an individual <code>exacct</code> accounting file object. An <code>exacct</code> object is represented as an opaque reference blessed into the appropriate <code>Sun::Solaris::Exacct::Object</code> subclass. This module is further subdivided into the object types <code>Item</code> and <code>Group</code> . At this level, there are methods to access the <code>ea_match_object_catalog(3EXACCT)</code> and <code>ea_attach_to_object(3EXACCT)</code> functions.	<code>Exacct::Object(3PERL)</code>
Sun::Solaris::Exacct::Object::Item	This module provides object-oriented methods to access an individual <code>exacct</code> accounting file <code>Item</code> . Objects of this type inherit from <code>Sun::Solaris::Exacct::Object</code> .	<code>Exacct::Object::Item(3PERL)</code>
Sun::Solaris::Exacct::Object::Group	This module provides object-oriented methods to access an individual <code>exacct</code> accounting file <code>Group</code> . Objects of this type inherit from <code>Sun::Solaris::Exacct::Object</code> . These objects provide access to the <code>ea_attach_to_group(3EXACCT)</code> function. The <code>Items</code> contained within the <code>Group</code> are presented as a Perl array.	<code>Exacct::Object::Group(3PERL)</code>
Sun::Solaris::Kstat	This module provides a Perl tied hash interface to the <code>kstat</code> facility. A usage example for this module can be found in <code>/bin/kstat</code> , which is written in Perl.	<code>Kstat(3PERL)</code>

For examples that show how to use the modules described in the previous table, see [“Using the Perl Interface to libexacct”](#) on page 67.

Administering Extended Accounting (Tasks)

This chapter describes how to administer the extended accounting subsystem. The following procedures are covered.

Administering the Extended Accounting Facility (Task Map)

Task	Description	For Instructions
Activate the extended accounting facility.	Use extended accounting to monitor resource consumption by each project running on your system. You can use the <i>extended accounting</i> subsystem to capture historical data for tasks, processes, and flows.	“How to Activate Extended Accounting for Processes, Tasks, and Flows” on page 64, “How to Activate Extended Accounting With a Startup Script” on page 64
Display extended accounting status.	Determine the status of the extended accounting facility.	“How to Display Extended Accounting Status” on page 65
View available accounting resources.	View the accounting resources available on your system.	“How to View Available Accounting Resources” on page 65
Deactivate the process, task, and flow accounting facility.	Turn off the extended accounting functionality.	“How to Deactivate Process, Task, and Flow Accounting” on page 66

Task	Description	For Instructions
Use the Perl interface to the extended accounting facility.	Use the Perl interface to to develop customized reporting and extraction scripts.	“Using the Perl Interface to libexacct” on page 67

Using Extended Accounting Functionality

▼ How to Activate Extended Accounting for Processes, Tasks, and Flows

To activate the extended accounting facility for tasks, processes, and flows, use the `acctadm` command. The optional final parameter to `acctadm` indicates whether the command should act on the process, system task, or flow accounting components of the extended accounting facility. See `acctadm(1M)` for more information.

1. Become superuser or assume an equivalent role.

Roles contain authorizations and privileged commands. For more information about roles, see “Using the Solaris Management Tools With RBAC (Task Map)” in *System Administration Guide: Basic Administration*.

2. Activate extended accounting for processes.

```
# acctadm -e extended -f /var/adm/exacct/proc process
```

3. Activate extended accounting for tasks.

```
# acctadm -e extended,mstate -f /var/adm/exacct/task task
```

4. Activate extended accounting for flows.

```
# acctadm -e extended -f /var/adm/exacct/flow flow
```

How to Activate Extended Accounting With a Startup Script

Activate extended accounting on an ongoing basis by linking the `/etc/init.d/acctadm` script into `/etc/rc2.d`.

```
# ln -s /etc/init.d/acctadm /etc/rc2.d/Snacctadm
# ln -s /etc/init.d/acctadm /etc/rc2.d/Knacctadm
```


The *n* variable is replaced by a number.

You must manually activate extended accounting at least once to set up the configuration.

See “Extended Accounting Configuration” on page 58 for information on accounting configuration.

How to Display Extended Accounting Status

Type `acctadm` without arguments to display the current status of the extended accounting facility.

```
# acctadm
      Task accounting: active
      Task accounting file: /var/adm/exacct/task
      Tracked task resources: extended
      Untracked task resources: none
      Process accounting: active
      Process accounting file: /var/adm/exacct/proc
      Tracked process resources: extended
      Untracked process resources: host
      Flow accounting: active
      Flow accounting file: /var/adm/exacct/flow
      Tracked flow resources: extended
      Untracked flow resources: none
```

In the previous example, system task accounting is active in extended mode and `mstate` mode. Process and flow accounting are active in extended mode.

Note – In the context of extended accounting, microstate (`mstate`) refers to the extended data, associated with microstate process transitions, that is available in the process usage file (see `proc(4)`). This data provides much more detail about the activities of the process than basic or extended records.

How to View Available Accounting Resources

Available resources can vary from system to system, and from platform to platform. Use the `acctadm` command with the `-r` option to view the accounting resources available on your system.

```
# acctadm -r
process:
extended pid,uid,gid,cpu,time,command,TTY,projid,taskid,ancpid,wait-status,zone,flag,
memory,mstate displays as one line
basic pid,uid,gid,cpu,time,command,TTY,flag
```

```

task:
extended taskid,projid,cpu,time,host,mstate,anctaskid,zone
basic    taskid,projid,cpu,time
flow:
extended saddr,daddr,sport,dport,proto,dsfield,nbytes,npkts,action,ctime,lseen,projid,uid
basic    saddr,daddr,sport,dport,proto,nbytes,npkts,action

```

▼ How to Deactivate Process, Task, and Flow Accounting

To deactivate process, task, and flow accounting, turn off each of them individually by using the `acctadm` command with the `-x` option.

1. Become superuser or assume an equivalent role.

Roles contain authorizations and privileged commands. For more information about roles, see “Using the Solaris Management Tools With RBAC (Task Map)” in *System Administration Guide: Basic Administration*.

2. Turn off process accounting.

```
# acctadm -x process
```

3. Turn off task accounting.

```
# acctadm -x task
```

4. Turn off flow accounting.

```
# acctadm -x flow
```

5. Verify that task accounting, process accounting, and flow accounting have been turned off.

```

# acctadm
    Task accounting: inactive
    Task accounting file: none
    Tracked task resources: extended
    Untracked task resources: none
    Process accounting: inactive
    Process accounting file: none
    Tracked process resources: extended
    Untracked process resources: host
    Flow accounting: inactive
    Flow accounting file: none
    Tracked flow resources: extended
    Untracked flow resources: none

```

Using the Perl Interface to libexacct

How to Recursively Print the Contents of an exact Object

Use the following code to recursively print the contents of an exact object. Note that this capability is provided by the library as the

Sun::Solaris::Exacct::Object::dump() function. This capability is also available through the ea_dump_object() convenience function.

```
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., an item or
        # group), then the value method will return a reference
        # to the appropriate sort of perl object
    }
}
```

```

# (Exacct::Object::Item or Exacct::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);
}
}

```

How to Create a New Group Record and Write It to a File

Use this script to create a new group record and write it to a file named /tmp/exacct.

```

#!/usr/perl5/5.6.1/bin/perl

use strict;
use warnings;

```

```

use Sun::Solaris::Exacct qw(:EXACCT_ALL);
# 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/rec" ],
);

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

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

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

# Open the exacct file, write the record & close.
my $f = ea_new_file('/tmp/exacct', &O_RDWR | &O_CREAT | &O_TRUNC)
    || die("create /tmp/exacct failed: ", ea_error_str(), "\n");
$f->write($group);
$f = undef;

```

How to Print the Contents of an exacct File

Use the following Perl script to print the contents of an exacct file.

```

#!/usr/perl5/5.6.1/bin/perl

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

die("Usage is dumpexacct <exacct file>\n") unless (@ARGV == 1);

# 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);

```

Example Output From Sun::Solaris::Exacct::Object->dump()

Here is example output produced by running
Sun::Solaris::Exacct::Object->dump() on the file created in [“How to Create a New Group Record and Write It to a File”](#) on page 68.

```

Creator: root
Hostname: localhost
GROUP
  Catalog = EXT_GROUP|EXC_DEFAULT|EXD_NONE
  ITEM
    Catalog = EXT_STRING|EXC_DEFAULT|EXD_CREATOR
    Value = me
  ITEM
    Catalog = EXT_UINT32|EXC_DEFAULT|EXD_PROC_PID
    Value = 845523
  ITEM
    Catalog = EXT_UINT32|EXC_DEFAULT|EXD_PROC_UID
    Value = 37845
  ITEM
    Catalog = EXT_UINT32|EXC_DEFAULT|EXD_PROC_GID
    Value = 10
  ITEM
    Catalog = EXT_STRING|EXC_DEFAULT|EXD_PROC_COMMAND
    Value = /bin/rec
ENDGROUP

```

Resource Controls (Overview)

After you determine the resource consumption of workloads on your system as described in [Chapter 4](#), you can place boundaries on resource usage. Boundaries prevent workloads from over-consuming resources. The *resource controls* facility is the constraint mechanism that is used for this purpose.

This chapter covers the following topics.

- “Resource Controls Concepts” on page 71
- “Configuring Resource Controls and Attributes” on page 73
- “Applying Resource Controls” on page 83
- “Temporarily Updating Resource Control Values on a Running System” on page 83
- “Commands Used With Resource Controls” on page 84

For information about how to administer resource controls, see [Chapter 7](#).

Resource Controls Concepts

In the Solaris Operating System, the concept of a per-process resource limit has been extended to the task and project entities described in [Chapter 2](#). These enhancements are provided by the resource controls (rctl) facility. In addition, allocations that were set through the `/etc/system` tunables are now automatic or configured through the resource controls mechanism as well.

A resource control is identified by the prefix `project`, `task`, or `process`. Resource controls can be observed on a system-wide basis. It is possible to update resource control values on a running system.

In a zones environment, the prefix `zone` identifies a zone-wide resource control. See “Resource Type Properties” on page 214 for information on available zone-wide resource controls.

For a list of the standard resource controls that are available in this release, see [“Available Resource Controls”](#) on page 74.

Resource Limits and Resource Controls

UNIX systems have traditionally provided a resource limit facility (*rlimit*). The *rlimit* facility allows administrators to set one or more numerical limits on the amount of resources a process can consume. These limits include per-process CPU time used, per-process core file size, and per-process maximum heap size. *Heap size* is the amount of scratch memory that is allocated for the process data segment.

The resource controls facility provides compatibility interfaces for the resource limits facility. Existing applications that use resource limits continue to run unchanged. These applications can be observed in the same way as applications that are modified to take advantage of the resource controls facility.

Interprocess Communication and Resource Controls

Processes can communicate with each other by using one of several types of interprocess communication (IPC). IPC allows information transfer or synchronization to occur between processes. Prior to the Solaris 10 release, IPC tunable parameters were set by adding an entry to the `/etc/system` file. The resource controls facility now provides resource controls that define the behavior of the kernel's IPC facilities. These resource controls replace the `/etc/system` tunables.

To observe which IPC objects are contributing to a project's usage, use the `ipcs` command with the `-J` option. See [“How to Use `ipcs`”](#) on page 92 to view an example display. For more information about the `ipcs` command, see `ipcs(1)`.

For information about Solaris system tuning, see the *Solaris Tunable Parameters Reference Manual*.

Resource Control Constraint Mechanisms

Resource controls provide a mechanism for the constraint of system resources. Processes, tasks, and projects can be prevented from consuming amounts of specified system resources. This mechanism leads to a more manageable system by preventing over-consumption of resources.

Constraint mechanisms can be used to support capacity-planning processes. An encountered constraint can provide information about application resource needs without necessarily denying the resource to the application.

Project Attribute Mechanisms

Resource controls can also serve as a simple attribute mechanism for resource management facilities. For example, the number of CPU shares made available to a project in the fair share scheduler (FSS) scheduling class is defined by the `project.cpu-shares` resource control. Because the project is assigned a fixed number of shares by the control, the various actions associated with exceeding a control are not relevant. In this context, the current value for the `project.cpu-shares` control is considered an attribute on the specified project.

Another type of project attribute is used to regulate the resource consumption of physical memory by collections of processes attached to a project. These attributes have the prefix `rcap`, for example, `rcap.max-rss`. Like a resource control, this type of attribute is configured in the `project` database. However, while resource controls are synchronously enforced by the kernel, resource caps are asynchronously enforced at the user level by the resource cap enforcement daemon, `rcapd`. For information on `rcapd`, see [Chapter 10](#) and `rcapd(1M)`.

The `project.pool` attribute is used to specify a pool binding for a project. For more information on resource pools, see [Chapter 12](#).

Configuring Resource Controls and Attributes

The resource controls facility is configured through the `project` database. See [Chapter 2](#). Resource controls and other attributes are set in the final field of the `project` database entry. The values associated with each resource control are enclosed in parentheses, and appear as plain text separated by commas. The values in parentheses constitute an “action clause.” Each action clause is composed of a privilege level, a threshold value, and an action that is associated with the particular threshold. Each resource control can have multiple action clauses, which are also separated by commas. The following entry defines a per-task lightweight process limit and a per-process maximum CPU time limit on a project entity. The `process.max-cpu-time` would send a process a `SIGTERM` after the process ran for 1 hour, and a `SIGKILL` if the process continued to run for a total of 1 hour and 1 minute. See [Table 6-2](#).

```
development:101:Developers:::task.max-lwps=(privileged,10,deny);
  process.max-cpu-time=(basic,3600,signal=TERM),(priv,3660,signal=KILL)
  typed as one line
```

Note – On systems that have zones enabled, zone-wide resource controls are specified in the zone configuration using a slightly different format. See [“Zone Configuration Data” on page 212](#) for more information.

The `rctladm` command allows you to make runtime interrogations of and modifications to the resource controls facility, with global scope. The `prctl` command allows you to make runtime interrogations of and modifications to the resource controls facility, with local scope.

For more information, see [“Global and Local Actions on Resource Control Values” on page 79](#), `rctladm(1M)` and `prctl(1)`.

Note – In a zones environment, you cannot use `rctladm` in a non-global zone to modify settings. You can use `rctladm` in a non-global zone to view the global logging state of each resource control.

Available Resource Controls

A list of the standard resource controls that are available in this release is shown in the following table.

The table describes the resource that is constrained by each control. The table also identifies the default units that are used by the `project` database for that resource. The default units are of two types:

- Quantities represent a limited amount.
- Indexes represent a maximum valid identifier.

Thus, `project.cpu-shares` specifies the number of shares to which the project is entitled. `process.max-file-descriptor` specifies the highest file number that can be assigned to a process by the `open(2)` system call.

TABLE 6–1 Standard Resource Controls

Control Name	Description	Default Unit
<code>project.cpu-shares</code>	Number of CPU shares granted to this project for use with the fair share scheduler (see FSS(7)).	Quantity (shares)

TABLE 6–1 Standard Resource Controls (Continued)

Control Name	Description	Default Unit
<code>project.max-crypto-memory</code>	Total amount of kernel memory that can be used by <code>libpkcs11</code> for hardware crypto acceleration. Allocations for kernel buffers and session-related structures are charged against this resource control.	Size (bytes)
<code>project.max-device-locked-memory</code>	Total amount of locked memory allowed.	Size (bytes)
<code>project.max-port-ids</code>	Maximum allowable number of event ports.	Quantity (number of event ports)
<code>project.max-shm-ids</code>	Maximum number of shared memory IDs allowed for this project.	Quantity (shared memory IDs)
<code>project.max-sem-ids</code>	Maximum number of semaphore IDs allowed for this project.	Quantity (semaphore IDs)
<code>project.max-msg-ids</code>	Maximum number of message queue IDs allowed for this project.	Quantity (message queue IDs)
<code>project.max-shm-memory</code>	Total amount of shared memory allowed for this project.	Size (bytes)
<code>project.max-lwps</code>	Maximum number of LWPs simultaneously available to this project.	Quantity (LWPs)
<code>project.max-tasks</code>	Maximum number of tasks allowable in this project.	Quantity (number of tasks)
<code>project.max-contracts</code>	Maximum number of contracts allowed in this project.	Quantity (contracts)
<code>task.max-cpu-time</code>	Maximum CPU time that is available to this task's processes.	Time (seconds)
<code>task.max-lwps</code>	Maximum number of LWPs simultaneously available to this task's processes.	Quantity (LWPs)
<code>process.max-cpu-time</code>	Maximum CPU time that is available to this process.	Time (seconds)

TABLE 6-1 Standard Resource Controls *(Continued)*

Control Name	Description	Default Unit
<code>process.max-file-descriptor</code>	Maximum file descriptor index available to this process.	Index (maximum file descriptor)
<code>process.max-file-size</code>	Maximum file offset available for writing by this process.	Size (bytes)
<code>process.max-core-size</code>	Maximum size of a core file created by this process.	Size (bytes)
<code>process.max-data-size</code>	Maximum heap memory available to this process.	Size (bytes)
<code>process.max-stack-size</code>	Maximum stack memory segment available to this process.	Size (bytes)
<code>process.max-address-space</code>	Maximum amount of address space, as summed over segment sizes, that is available to this process.	Size (bytes)
<code>process.max-port-events</code>	Maximum allowable number of events per event port.	Quantity (number of events)
<code>process.max-sem-nsems</code>	Maximum number of semaphores allowed per semaphore set.	Quantity (semaphores per set)
<code>process.max-sem-ops</code>	Maximum number of semaphore operations allowed per <code>semop</code> call (value copied from the resource control at <code>semget ()</code> time).	Quantity (number of operations)
<code>process.max-msg-qbytes</code>	Maximum number of bytes of messages on a message queue (value copied from the resource control at <code>msgget ()</code> time).	Size (bytes)
<code>process.max-msg-messages</code>	Maximum number of messages on a message queue (value copied from the resource control at <code>msgget ()</code> time).	Quantity (number of messages)

The following zone-wide resource controls are available in a zones environment:

- `zone.cpu-shares`
- `zone.max-lwps`

For information on zone-wide resource controls, see “Resource Type Properties” on page 214.

Units Support

Global flags that identify resource control types are defined for all resource controls. The flags are used by the system to communicate basic type information to applications such as the `prctl` command. Applications use the information to determine the following:

- The unit strings that are appropriate for each resource control
- The correct scale to use when interpreting scaled values

The following global flags are available:

Global Flag	Resource Control Type String	Modifier	Scale
RCTL_GLOBAL_BYTES	bytes	B	1
		KB	2^{10}
		MB	2^{20}
		GB	2^{30}
		TB	2^{40}
		PB	2^{50}
		EB	2^{60}
RCTL_GLOBAL_SECONDS	seconds	s	1
		Ks	10^3
		Ms	10^6
		Gs	10^9
		Ts	10^{12}
		Ps	10^{15}
		Es	10^{18}

Global Flag	Resource Control Type String	Modifier	Scale
RCTL_GLOBAL_COUNT	count	none	1
		K	10 ³
		M	10 ⁶
		G	10 ⁹
		T	10 ¹²
		P	10 ¹⁵
		E	10 ¹⁸

Scaled values can be used with resource controls. The following example shows a scaled threshold value:

```
task.max-lwps=(priv,1K,deny)
```

Note – Unit modifiers are accepted by the `prctl`, `projadd`, and `projmod` commands. You cannot use unit modifiers in the `project` database itself.

Resource Control Values and Privilege Levels

A threshold value on a resource control constitutes an enforcement point where local actions can be triggered or global actions, such as logging, can occur.

Each threshold value on a resource control must be associated with a privilege level. The privilege level must be one of the following three types.

- Basic, which can be modified by the owner of the calling process
- Privileged, which can be modified only by privileged (superuser) callers
- System, which is fixed for the duration of the operating system instance

A resource control is guaranteed to have one system value, which is defined by the system, or resource provider. The system value represents how much of the resource the current implementation of the operating system is capable of providing.

Any number of privileged values can be defined, and only one basic value is allowed. Operations that are performed without specifying a privilege value are assigned a basic privilege by default.

The privilege level for a resource control value is defined in the privilege field of the resource control block as `RCTL_BASIC`, `RCTL_PRIVILEGED`, or `RCTL_SYSTEM`. See `setrctl(2)` for more information. You can use the `prctl` command to modify values that are associated with basic and privileged levels.

Global and Local Actions on Resource Control Values

There are two categories of actions on resource control values: global and local.

Global Actions on Resource Control Values

Global actions apply to resource control values for every resource control on the system. You can use the `rctladm` command described in the `rctladm(1M)` man page to perform the following actions:

- Display the global state of active system resource controls
- Set global logging actions

You can disable or enable the global logging action on resource controls. You can set the `syslog` action to a specific degree by assigning a severity level, `syslog=level`. The possible settings for `level` are as follows:

- `debug`
- `info`
- `notice`
- `warning`
- `err`
- `crit`
- `alert`
- `emerg`

By default, there is no global logging of resource control violations.

Local Actions on Resource Control Values

Local actions are taken on a process that attempts to exceed the control value. For each threshold value that is placed on a resource control, you can associate one or more actions. There are three types of local actions: `none`, `deny`, and `signal=`. These three actions are used as follows:

<code>none</code>	No action is taken on resource requests for an amount that is greater than the threshold. This action is useful for monitoring resource usage without affecting the progress of applications. You can also enable a global message that displays when the resource control is exceeded, although the process exceeding the threshold is not affected.
<code>deny</code>	You can deny resource requests for an amount that is greater than the threshold. For example, a <code>task.max-lwps</code> resource control with action <code>deny</code> causes a <code>fork</code> system call to fail if the new process would exceed the control value. See the <code>fork(2)</code> man page.
<code>signal=</code>	You can enable a global signal message action when the resource control is exceeded. A signal is sent to the process when the threshold value is

exceeded. Additional signals are not sent if the process consumes additional resources. Available signals are listed in [Table 6–2](#).

Not all of the actions can be applied to every resource control. For example, a process cannot exceed the number of CPU shares assigned to the project of which it is a member. Therefore, a deny action is not allowed on the `project.cpu-shares` resource control.

Due to implementation restrictions, the global properties of each control can restrict the range of available actions that can be set on the threshold value. (See the `rctladm(1M)` man page.) A list of available signal actions is presented in the following table. For additional information about signals, see the `signal(3HEAD)` man page.

TABLE 6–2 Signals Available to Resource Control Values

Signal	Description	Notes
SIGABRT	Terminate the process.	
SIGHUP	Send a hangup signal. Occurs when carrier drops on an open line. Signal sent to the process group that controls the terminal.	
SIGTERM	Terminate the process. Termination signal sent by software.	
SIGKILL	Terminate the process and kill the program.	
SIGSTOP	Stop the process. Job control signal.	
SIGXRES	Resource control limit exceeded. Generated by resource control facility.	
SIGXFSZ	Terminate the process. File size limit exceeded.	Available only to resource controls with the <code>RCTL_GLOBAL_FILE_SIZE</code> property (<code>process.max-file-size</code>). See <code>rctlblk_set_value(3C)</code> for more information.

TABLE 6-2 Signals Available to Resource Control Values *(Continued)*

Signal	Description	Notes
SIGXCPU	Terminate the process. CPU time limit exceeded.	Available only to resource controls with the RCTL_GLOBAL_CPU_TIME property (process.max-cpu-time). See rctlblk_set_value(3C) for more information.

Resource Control Flags and Properties

Each resource control on the system has a certain set of associated properties. This set of properties is defined as a set of flags, which are associated with all controlled instances of that resource. Global flags cannot be modified, but the flags can be retrieved by using either `rctladm` or the `getrctl` system call.

Local flags define the default behavior and configuration for a specific threshold value of that resource control on a specific process or process collective. The local flags for one threshold value do not affect the behavior of other defined threshold values for the same resource control. However, the global flags affect the behavior for every value associated with a particular control. Local flags can be modified, within the constraints supplied by their corresponding global flags, by the `prctl` command or the `setrctl` system call. See `setrctl(2)`.

For the complete list of local flags, global flags, and their definitions, see `rctlblk_set_value(3C)`.

To determine system behavior when a threshold value for a particular resource control is reached, use `rctladm` to display the global flags for the resource control. For example, to display the values for `process.max-cpu-time`, type the following:

```
$ rctladm process.max-cpu-time
process.max-cpu-time syslog=off [ lowerable no-deny cpu-time inf seconds ]
```

The global flags indicate the following.

- `lowerable` Superuser privileges are not required to lower the privileged values for this control.
- `no-deny` Even when threshold values are exceeded, access to the resource is never denied.
- `cpu-time` SIGXCPU is available to be sent when threshold values of this resource are reached.
- `seconds` The time value for the resource control.

Use the `prctl` command to display local values and actions for the resource control.

```

$ prctl -n process.max-cpu-time $$
process 353939: -ksh
NAME      PRIVILEGE  VALUE   FLAG   ACTION          RECIPIENT
process.max-cpu-time
  privileged 18.4Es   inf    signal=XCPU     -
  system    18.4Es   inf    none

```

The max (RCTL_LOCAL_MAXIMAL) flag is set for both threshold values, and the inf (RCTL_GLOBAL_INFINITE) flag is defined for this resource control. An inf value has an infinite quantity. The value is never enforced. Hence, as configured, both threshold quantities represent infinite values that are never exceeded.

Resource Control Enforcement

More than one resource control can exist on a resource. A resource control can exist at each containment level in the process model. If resource controls are active on the same resource at different container levels, the smallest container's control is enforced first. Thus, action is taken on `process.max-cpu-time` before `task.max-cpu-time` if both controls are encountered simultaneously.

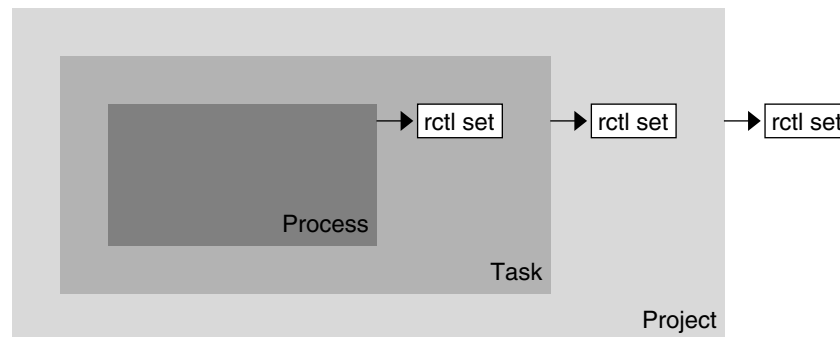


FIGURE 6-1 Process Collectives, Container Relationships, and Their Resource Control Sets

Global Monitoring of Resource Control Events

Often, the resource consumption of processes is unknown. To get more information, try using the global resource control actions that are available with the `rctladm` command. Use `rctladm` to establish a `syslog` action on a resource control. Then, if any entity managed by that resource control encounters a threshold value, a system message is logged at the configured logging level. See [Chapter 7](#) and the `rctladm(1M)` man page for more information.

Applying Resource Controls

Each resource control listed in [Table 6-1](#) can be assigned to a project at login or when `newtask`, `su`, or the other project-aware launchers `at`, `batch`, or `cron` are invoked. Each command that is initiated is launched in a separate task with the invoking user's default project. See the man pages `login(1)`, `newtask(1)`, `at(1)`, `cron(1M)`, and `su(1M)` for more information.

Updates to entries in the `project` database, whether to the `/etc/project` file or to a representation of the database in a network name service, are not applied to currently active projects. The updates are applied when a new task joins the project through `login` or `newtask`.

Temporarily Updating Resource Control Values on a Running System

Values changed in the `project` database only become effective for new tasks that are started in a project. However, you can use the `rctladm` and `prctl` commands to update resource controls on a running system.

Updating Logging Status

The `rctladm` command affects the global logging state of each resource control on a system-wide basis. This command can be used to view the global state and to set up the level of `syslog` logging when controls are exceeded.

Updating Resource Controls

You can view and temporarily alter resource control values and actions on a per-process, per-task, or per-project basis by using the `prctl` command. A project, task, or process ID is given as input, and the command operates on the resource control at the level where the control is defined.

Any modifications to values and actions take effect immediately. However, these modifications apply to the current process, task, or project only. The changes are not recorded in the `project` database. If the system is restarted, the modifications are lost. Permanent changes to resource controls must be made in the `project` database.

All resource control settings that can be modified in the `project` database can also be modified with the `prctl` command. Both basic and privileged values can be added or be deleted. Their actions can also be modified. By default, the basic type is assumed for all set operations, but processes and users with superuser privileges can also modify privileged resource controls. System resource controls cannot be altered.

Commands Used With Resource Controls

The commands that are used with resource controls are shown in the following table.

Command Reference	Description
<code>ipcs(1)</code>	Allows you to observe which IPC objects are contributing to a project's usage
<code>prctl(1)</code>	Allows you to make runtime interrogations of and modifications to the resource controls facility, with local scope
<code>rctladm(1M)</code>	Allows you to make runtime interrogations of and modifications to the resource controls facility, with global scope

Administering Resource Controls (Tasks)

This chapter describes how to administer the resource controls facility.

For an overview of the resource controls facility, see [Chapter 6](#).

Administering Resource Controls (Task Map)

Task	Description	For Instructions
Set resource controls.	Set resource controls for a project in the <code>/etc/project</code> file.	“Setting Resource Controls” on page 86
Get or revise the resource control values for active processes, tasks, or projects, with local scope.	Make runtime interrogations of and modifications to the resource controls associated with an active process, task, or project on the system.	“Using the <code>prctl</code> Command” on page 89
On a running system, view or update the global state of resource controls.	View the global logging state of each resource control on a system-wide basis. Also set up the level of <code>syslog</code> logging when controls are exceeded.	“Using <code>rctladm</code>” on page 91

Task	Description	For Instructions
Report status of active interprocess communication (IPC) facilities.	Display information about active interprocess communication (IPC) facilities. Observe which IPC objects are contributing to a project's usage.	"Using ipcs" on page 92
Determine whether a web server is allocated sufficient CPU capacity.	Set a global action on a resource control. This action enables you to receive notice of any entity that has a resource control value that is set too low.	"How to Determine Whether a Web Server Is Allocated Enough CPU Capacity" on page 93

Setting Resource Controls

▼ How to Set the Maximum Number of LWPs for Each Task in a Project

This procedure adds a project named `x-files` to the `/etc/project` file and sets a maximum number of LWPs for a task created in the project.

1. Become superuser or assume an equivalent role.

Roles contain authorizations and privileged commands. For more information about roles, see *"Using the Solaris Management Tools With RBAC (Task Map)"* in *System Administration Guide: Basic Administration*.

2. Use the `projadd` command with the `-K` option to create a project called `x-files`. Set the maximum number of LWPs for each task created in the project to 3.

```
# projadd -K 'task.max-lwps=(privileged,3,deny)' x-files
```

3. View the entry in the `/etc/project` file by using one of the following methods:

■ Type:

```
# projects -l
system
    projid : 0
    comment: ""
```

```

        users : (none)
        groups : (none)
        attribs:
    .
    .
    .
x-files
    projid : 100
    comment: ""
    users  : (none)
    groups : (none)
    attribs: task.max-lwps=(privileged,3,deny)

```

■ Type:

```

# cat etc/project
system:0:System:::
.
.
.
x-files:100:::task.max-lwps=(privileged,3,deny)

```

When superuser creates a new task in project `x-files` by joining the project with `newtask`, superuser will not be able to create more than three LWPs while running in this task. This is shown in the following annotated sample session.

```

# newtask -p x-files csh

# prctl -n task.max-lwps $$
process: 111107: csh
NAME      PRIVILEGE  VALUE   FLAG   ACTION           RECIPIENT
task.max-lwps
          privileged    3      -     deny            -
          system      2.15G   max   deny            -

# id -p
uid=0(root) gid=1(other) projid=100(x-files)

# ps -o project,taskid -p $$
PROJECT TASKID
x-files   73

# csh          /* creates second LWP */

# csh          /* creates third LWP */

# csh          /* cannot create more LWPs */
Vfork failed
#

```

▼ How to Set Multiple Controls on a Project

The `/etc/project` file can contain settings for multiple resource controls for each project as well as multiple threshold values for each control. Threshold values are defined in action clauses, which are comma-separated for multiple values.

1. **Become superuser or assume an equivalent role.**

Roles contain authorizations and privileged commands. For more information about roles, see “Using the Solaris Management Tools With RBAC (Task Map)” in *System Administration Guide: Basic Administration*.

2. **Use the `projmod` command with the `-s` and `-K` options to set resource controls on project `x-files`:**

```
# projmod -s -K 'task.max-lwps=(basic,10,none), (privileged,500,deny);
process.max-file-descriptor=(basic,128,deny)' x-files one line in file
```

The following controls are set:

- A `basic` control with no action on the maximum LWPs per task.
- A privileged deny control on the maximum LWPs per task. This control causes any LWP creation that exceeds the maximum to fail, as shown in the previous example “How to Set the Maximum Number of LWPs for Each Task in a Project” on page 86.
- A limit on the maximum file descriptors per process at the `basic` level, which forces the failure of any open call that exceeds the maximum.

3. **View the entry in the file by using one of the following methods:**

- **Type:**

```
# projects -l
.
.
.
x-files
    projid : 100
    comment: ""
    users  : (none)
    groups : (none)
    attribs: process.max-file-descriptor=(basic,128,deny)
            task.max-lwps=(basic,10,none), (privileged,500,deny) one
line in file
```

- **Type:**

```
# cat etc/project
.
.
.
x-files:100:::process.max-file-descriptor=(basic,128,deny);
task.max-lwps=(basic,10,none), (privileged,500,deny) one line in file
```

Using the `prctl` Command

Use the `prctl` command to make runtime interrogations of and modifications to the resource controls associated with an active process, task, or project on the system. See the `prctl(1)` man page for more information.

How to Use `prctl` to Display Information

The following command displays the maximum file descriptor for the current shell that is running.

```
# prctl -n process.max-file-descriptor $$
process: 110453: -sh
NAME      PRIVILEGE      VALUE      FLAG      ACTION      RECIPIENT
process.max-file-descriptor
  basic           256         -         deny      110453
  privileged      65.5K       -         deny      -
  system          2.15G      max        deny
```

How to Use `prctl` to Temporarily Change a Value

This example procedure uses the `prctl` command to temporarily add a new privileged value to deny the use of more than three LWPs per project for the `x-files` project. The result is comparable to the result in [“How to Set the Maximum Number of LWPs for Each Task in a Project”](#) on page 86.

1. Become superuser or assume an equivalent role.

Roles contain authorizations and privileged commands. For more information about roles, see *“Using the Solaris Management Tools With RBAC (Task Map)”* in *System Administration Guide: Basic Administration*.

2. Use `newtask` to join the `x-files` project.

```
# newtask -p x-files
```

3. Use the `id` command with the `-p` option to verify that the correct project has been joined.

```
# id -p
uid=0(root) gid=1(other) projid=101(x-files)
```

4. Add a new privileged value for `project.max-lwps` that limits the number of LWPs to three.

```
# prctl -n project.max-lwps -t privileged -v 3 -e deny -i project x-files
```

5. Verify the result.

```
# prctl -n project.max-lwps -i project x-files
process: 111108: csh
NAME      PRIVILEGE      VALUE      FLAG      ACTION      RECIPIENT
project.max-lwps
  privileged          3          -      deny      -
  system             2.15G      max      deny      -
```

How to Use `prctl` to Lower a Resource Control Value

1. Become superuser or assume an equivalent role.

Roles contain authorizations and privileged commands. For more information about roles, see “Using the Solaris Management Tools With RBAC (Task Map)” in *System Administration Guide: Basic Administration*.

2. Use the `prctl` command with the `-r` option to change the lowest value of the `process.max-file-descriptor` resource control.

```
# prctl -n process.max-file-descriptor -r -v 128 $$
```

▼ How to Use `prctl` to Display, Replace, and Verify the Value of a Control on a Project

1. Become superuser or assume an equivalent role.

Roles contain authorizations and privileged commands. For more information about roles, see “Using the Solaris Management Tools With RBAC (Task Map)” in *System Administration Guide: Basic Administration*.

2. Display the value of `project.cpu-shares` in the project `group.staff`.

```
# prctl -n project.cpu-shares -i project group.staff
project: 2: group.staff
NAME      PRIVILEGE      VALUE      FLAG      ACTION      RECIPIENT
project.cpu-shares
  privileged          1          -      none      -
  system             65.5K      max      none
```

3. Replace the current `project.cpu-shares` value 1 with the value 10.

```
# prctl -n project.cpu-shares -v 10 -r -i project group.staff
```

4. Display the value of `project.cpu-shares` in the project `group.staff`.

```
# prctl -n project.cpu-shares -i project group.staff
project: 2: group.staff
NAME      PRIVILEGE      VALUE      FLAG      ACTION      RECIPIENT
```

```

project.cpu-shares
    privileged      10      -      none      -
    system          65.5K    max    none

```

Using rctladm

How to Use rctladm

Use the `rctladm` command to make runtime interrogations of and modifications to the global state of the resource controls facility. See the `rctladm(1M)` man page for more information.

For example, you can use `rctladm` with the `-e` option to enable the global `syslog` attribute of a resource control. When the control is exceeded, notification is logged at the specified `syslog` level. To enable the global `syslog` attribute of `process.max-file-descriptor`, type the following:

```
# rctladm -e syslog process.max-file-descriptor
```

When used without arguments, the `rctladm` command displays the global flags, including the global type flag, for each resource control.

```

# rctladm
process.max-port-events      syslog=off [ deny count ]
process.max-msg-messages     syslog=off [ deny count ]
process.max-msg-qbytes       syslog=off [ deny bytes ]
process.max-sem-ops          syslog=off [ deny count ]
process.max-sem-nsems        syslog=off [ deny count ]
process.max-address-space     syslog=off [ lowerable deny no-signal bytes ]
process.max-file-descriptor   syslog=off [ lowerable deny count ]
process.max-core-size         syslog=off [ lowerable deny no-signal bytes ]
process.max-stack-size        syslog=off [ lowerable deny no-signal bytes ]
.
.
.

```

Using `ipcs`

How to Use `ipcs`

Use the `ipcs` utility to display information about active interprocess communication (IPC) facilities. See the `ipcs(1)` man page for more information.

You can use `ipcs` with the `-J` option to see which project's limit an IPC object is allocated against.

```
# ipcs -J
IPC status from <running system> as of Wed Mar 26 18:53:15 PDT 2003
T          ID      KEY      MODE      OWNER     GROUP     PROJECT
Message Queues:
Shared Memory:
m          3600     0      --rw-rw-rw-  uname    staff    x-files
m           201     0      --rw-rw-rw-  uname    staff    x-files
m          1802     0      --rw-rw-rw-  uname    staff    x-files
m           503     0      --rw-rw-rw-  uname    staff    x-files
m           304     0      --rw-rw-rw-  uname    staff    x-files
m           605     0      --rw-rw-rw-  uname    staff    x-files
m            6     0      --rw-rw-rw-  uname    staff    x-files
m          107     0      --rw-rw-rw-  uname    staff    x-files
Semaphores:
s            0     0      --rw-rw-rw-  uname    staff    x-files
```

Capacity Warnings

A global action on a resource control enables you to receive notice of any entity that is tripping over a resource control value that is set too low.

For example, assume you want to determine whether a web server possesses sufficient CPUs for its typical workload. You could analyze `sar` data for idle CPU time and load average. You could also examine extended accounting data to determine the number of simultaneous processes that are running for the web server process.

However, an easier approach is to place the web server in a task. You can then set a global action, using `syslog`, to notify you whenever a task exceeds a scheduled number of LWPs appropriate for the machine's capabilities.

See the `sar(1)` man page for more information.

▼ How to Determine Whether a Web Server Is Allocated Enough CPU Capacity

1. Use the `prctl` command to place a privileged (superuser-owned) resource control on the tasks that contain an `httpd` process. Limit each task's total number of LWPs to 40, and disable all local actions.

```
# prctl -n task.max-lwps -v 40 -t privileged -d all `pgrep httpd`
```

2. Enable a system log global action on the `task.max-lwps` resource control.

```
# rctladm -e syslog task.max-lwps
```

3. Observe whether the workload trips the resource control.

If it does, you will see `/var/adm/messages` such as:

```
Jan  8 10:15:15 testmachine unix: [ID 859581 kern.notice]  
NOTICE: privileged rctl task.max-lwps exceeded by task 19
```


Fair Share Scheduler (Overview)

The analysis of workload data can indicate that a particular workload or group of workloads is monopolizing CPU resources. If these workloads are not violating resource constraints on CPU usage, you can modify the allocation policy for CPU time on the system. The fair share scheduling class described in this chapter enables you to allocate CPU time based on shares instead of the priority scheme of the timesharing (TS) scheduling class.

This chapter covers the following topics.

- “Introduction to the Scheduler” on page 95
- “CPU Share Definition” on page 96
- “CPU Shares and Process State” on page 97
- “CPU Share Versus Utilization” on page 97
- “CPU Share Examples” on page 98
- “FSS Configuration” on page 100
- “FSS and Processor Sets” on page 102
- “Combining FSS With Other Scheduling Classes” on page 104
- “Setting the Scheduling Class for the System” on page 105
- “Scheduling Class in a Zones Environment” on page 105
- “Commands Used With FSS” on page 106

To begin using the fair share scheduler, see [Chapter 9](#).

Introduction to the Scheduler

A fundamental job of the operating system is to arbitrate which processes get access to the system’s resources. The process scheduler, which is also called the dispatcher, is the portion of the kernel that controls allocation of the CPU to processes. The scheduler supports the concept of scheduling classes. Each class defines a scheduling policy that is used to schedule processes within the class. The default scheduler in the

Solaris Operating System, the TS scheduler, tries to give every process relatively equal access to the available CPUs. However, you might want to specify that certain processes be given more resources than others.

You can use the *fair share scheduler* (FSS) to control the allocation of available CPU resources among workloads, based on their importance. This importance is expressed by the number of *shares* of CPU resources that you assign to each workload.

You give each project CPU shares to control the project's entitlement to CPU resources. The FSS guarantees a fair dispersion of CPU resources among projects that is based on allocated shares, independent of the number of processes that are attached to a project. The FSS achieves fairness by reducing a project's entitlement for heavy CPU usage and increasing its entitlement for light usage, in accordance with other projects.

The FSS consists of a kernel scheduling class module and class-specific versions of the `dispadm(1M)` and `pricntl(1)` commands. Project shares used by the FSS are specified through the `project.cpu-shares` property in the `project(4)` database.

Note – If you are using the `project.cpu-shares` resource control in a zones environment, see [“Zone Configuration Data” on page 212](#), [“Resource Controls Used in Non-Global Zones” on page 305](#), and [“Using the Fair Share Scheduler in a Zones Environment” on page 323](#).

CPU Share Definition

The term “share” is used to define a portion of the system's CPU resources that is allocated to a project. If you assign a greater number of CPU shares to a project, relative to other projects, the project receives more CPU resources from the fair share scheduler.

CPU shares are not equivalent to percentages of CPU resources. Shares are used to define the relative importance of workloads in relation to other workloads. When you assign CPU shares to a project, your primary concern is not the number of shares the project has. Knowing how many shares the project has in comparison with other projects is more important. You must also take into account how many of those other projects will be competing with it for CPU resources.

Note – Processes in projects with zero shares always run at the lowest system priority (0). These processes only run when projects with nonzero shares are not using CPU resources.

CPU Shares and Process State

In the Solaris system, a project workload usually consists of more than one process. From the fair share scheduler perspective, each project workload can be in either an *idle* state or an *active* state. A project is considered idle if none of its processes are using any CPU resources. This usually means that such processes are either *sleeping* (waiting for I/O completion) or stopped. A project is considered active if at least one of its processes is using CPU resources. The sum of shares of all active projects is used in calculating the portion of CPU resources to be assigned to projects.

The following formula shows how the FSS scheduler calculates per-project allocation of CPU resources.

$$\text{allocation}_{\text{project } i} = \frac{\text{shares}_{\text{project } i}}{\sum_{j=1 \dots n} (\text{shares}_{\text{project } j})}$$

j is the index among all active projects

FIGURE 8-1 FSS Scheduler Share Calculation

When more projects become active, each project's CPU allocation is reduced, but the proportion between the allocations of different projects does not change.

CPU Share Versus Utilization

Share allocation is not the same as utilization. A project that is allocated 50 percent of the CPU resources might average only a 20 percent CPU use. Moreover, shares serve to limit CPU usage only when there is competition from other projects. Regardless of how low a project's allocation is, it always receives 100 percent of the processing power if it is running alone on the system. Available CPU cycles are never wasted. They are distributed between projects.

The allocation of a small share to a busy workload might slow its performance. However, the workload is not prevented from completing its work if the system is not overloaded.

CPU Share Examples

Assume you have a system with two CPUs running two parallel CPU-bound workloads called *A* and *B*, respectively. Each workload is running as a separate project. The projects have been configured so that project *A* is assigned S_A shares, and project *B* is assigned S_B shares.

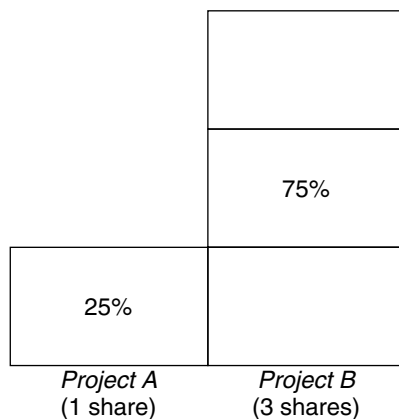
On average, under the traditional TS scheduler, each of the workloads that is running on the system would be given the same amount of CPU resources. Each workload would get 50 percent of the system's capacity.

When run under the control of the FSS scheduler with $S_A=S_B$, these projects are also given approximately the same amounts of CPU resources. However, if the projects are given different numbers of shares, their CPU resource allocations are different.

The next three examples illustrate how shares work in different configurations. These examples show that shares are only mathematically accurate for representing the usage if demand meets or exceeds available resources.

Example 1: Two CPU-Bound Processes in Each Project

If *A* and *B* each have two CPU-bound processes, and $S_A = 1$ and $S_B = 3$, then the total number of shares is $1 + 3 = 4$. In this configuration, given sufficient CPU demand, projects *A* and *B* are allocated 25 percent and 75 percent of CPU resources, respectively.



Example 2: No Competition Between Projects

If A and B have only *one* CPU-bound process each, and $S_A = 1$ and $S_B = 100$, then the total number of shares is 101. Each project cannot use more than one CPU because each project has only one running process. Because no competition exists between projects for CPU resources in this configuration, projects A and B are each allocated 50 percent of all CPU resources. In this configuration, CPU share values are irrelevant. The projects' allocations would be the same (50/50), even if both projects were assigned zero shares.

50%	50%
(1st CPU)	(2nd CPU)
<i>Project A</i> (1 share)	<i>Project B</i> (100 shares)

Example 3: One Project Unable to Run

If A and B have two CPU-bound processes each, and project A is given 1 share and project B is given 0 shares, then project B is not allocated any CPU resources and project A is allocated all CPU resources. Processes in B always run at system priority 0, so they never be able to run because processes in project A always have higher priorities.



FSS Configuration

Projects and Users

Projects are the workload containers in the FSS scheduler. Groups of users who are assigned to a project are treated as single controllable blocks. Note that you can create a project with its own number of shares for an individual user.

Users can be members of multiple projects that have different numbers of shares assigned. By moving processes from one project to another project, processes can be assigned CPU resources in varying amounts.

For more information on the `project(4)` database and name services, see [“project Database”](#) on page 35.

CPU Shares Configuration

The configuration of CPU shares is managed by the name service as a property of the project database.

When the first task (or process) that is associated with a project is created through the `setproject(3PROJECT)` library function, the number of CPU shares defined as resource control `project.cpu-shares` in the `project` database is passed to the kernel. A project that does not have the `project.cpu-shares` resource control defined is assigned one share.

In the following example, this entry in the `/etc/project` file sets the number of shares for project *x-files* to 5:

```
x-files:100:::project.cpu-shares=(privileged,5,none)
```

If you alter the number of CPU shares allocated to a project in the database when processes are already running, the number of shares for that project will not be modified at that point. The project must be restarted for the change to become effective.

If you want to temporarily change the number of shares assigned to a project without altering the project's attributes in the `project` database, use the `prctl` command. For example, to change the value of project *x-files*'s `project.cpu-shares` resource control to 3 while processes associated with that project are running, type the following:

```
# prctl -r -n project.cpu-shares -v 3 -i project x-files
See the prctl(1) man page for more information.
```

<code>-r</code>	Replaces the current value for the named resource control.
<code>-n name</code>	Specifies the name of the resource control.
<code>-v val</code>	Specifies the value for the resource control.
<code>-i idtype</code>	Specifies the ID type of the next argument.
<i>x-files</i>	Specifies the object of the change. In this instance, project <i>x-files</i> is the object.

Project `system` with project ID 0 includes all system daemons that are started by the boot-time initialization scripts. `system` can be viewed as a project with an unlimited number of shares. This means that `system` is always scheduled first, regardless of how many shares have been given to other projects. If you do not want the `system` project to have unlimited shares, you can specify a number of shares for this project in the `project` database.

As stated previously, processes that belong to projects with zero shares are always given zero system priority. Projects with one or more shares are running with priorities one and higher. Thus, projects with zero shares are only scheduled when CPU resources are available that are not requested by a nonzero share project.

The maximum number of shares that can be assigned to one project is 65535.

FSS and Processor Sets

The FSS can be used in conjunction with processor sets to provide more fine-grained controls over allocations of CPU resources among projects that run on each processor set than would be available with processor sets alone. The FSS scheduler treats processor sets as entirely independent partitions, with each processor set controlled independently with respect to CPU allocations.

The CPU allocations of projects running in one processor set are not affected by the CPU shares or activity of projects running in another processor set because the projects are not competing for the same resources. Projects only compete with each other if they are running within the same processor set.

The number of shares allocated to a project is system wide. Regardless of which processor set it is running on, each portion of a project is given the same amount of shares.

When processor sets are used, project CPU allocations are calculated for active projects that run within each processor set, as shown in the following figure.

$$\text{allocation}_{\text{project } x}^i = \frac{\text{shares}_{\text{project } x}^i}{\sum_{j=1 \dots n} (\text{shares}_{\text{project } j})}$$

j is the index among all active projects that run on processor set X

FIGURE 8-2 FSS Scheduler Share Calculation With Processor Sets

Project partitions that run on different processor sets might have different CPU allocations. The CPU allocation for each project partition in a processor set depends only on the allocations of other projects that run on the same processor set.

The performance and availability of applications that run within the boundaries of their processor sets are not affected by the introduction of new processor sets. The applications are also not affected by changes that are made to the share allocations of projects that run on other processor sets.

Empty processor sets (sets without processors in them) or processor sets without processes bound to them do not have any impact on the FSS scheduler behavior.

FSS and Processor Sets Examples

Assume that a server with eight CPUs is running several CPU-bound applications in projects *A*, *B*, and *C*. Project *A* is allocated one share, project *B* is allocated two shares, and project *C* is allocated three shares.

Project *A* is running only on processor set 1. Project *B* is running on processor sets 1 and 2. Project *C* is running on processor sets 1, 2, and 3. Assume that each project has enough processes to utilize all available CPU power. Thus, there is always competition for CPU resources on each processor set.

Project A 16.66% (1/6)	Project B 40% (2/5)	Project C 100% (3/3)
Project B 33.33% (2/6)		
Project C 50% (3/6)	Project C 60% (3/5)	
Processor Set #1 2 CPUs 25% of the system	Processor Set #2 4 CPUs 50% of the system	Processor Set #3 2 CPUs 25% of the system

The total system-wide project CPU allocations on such a system are shown in the following table.

Project	Allocation
Project A	$4\% = (1/6 \times 2/8)_{\text{pset1}}$
Project B	$28\% = (2/6 \times 2/8)_{\text{pset1}} + (2/5 \times 4/8)_{\text{pset2}}$
Project C	$67\% = (3/6 \times 2/8)_{\text{pset1}} + (3/5 \times 4/8)_{\text{pset2}} + (3/3 \times 2/8)_{\text{pset3}}$

These percentages do not match the corresponding amounts of CPU shares that are given to projects. However, within each processor set, the per-project CPU allocation ratios are proportional to their respective shares.

On the same system *without* processor sets, the distribution of CPU resources would be different, as shown in the following table.

Project	Allocation
Project A	16.66% = (1/6)
Project B	33.33% = (2/6)
Project C	50% = (3/6)

Combining FSS With Other Scheduling Classes

By default, the FSS scheduling class uses the same range of priorities (0 to 59) as the timesharing (TS), interactive (IA), and fixed priority (FX) scheduling classes. Therefore, you should avoid having processes from these scheduling classes share *the same* processor set. A mix of processes in the FSS, TS, IA, and FX classes could result in unexpected scheduling behavior.

With the use of processor sets, you can mix TS, IA, and FX with FSS in one system. However, all the processes that run on each processor set must be in *one* scheduling class, so they do not compete for the same CPUs. The FX scheduler in particular should not be used in conjunction with the FSS scheduling class unless processor sets are used. This action prevents applications in the FX class from using priorities high enough to starve applications in the FSS class.

You can mix processes in the TS and IA classes in the same processor set, or on the same system without processor sets.

The Solaris system also offers a real-time (RT) scheduler to users with superuser privileges. By default, the RT scheduling class uses system priorities in a different range (usually from 100 to 159) than FSS. Because RT and FSS are using *disjoint*, or non-overlapping, ranges of priorities, FSS can coexist with the RT scheduling class within the same processor set. However, the FSS scheduling class does not have any control over processes that run in the RT class.

For example, on a four-processor system, a single-threaded RT process can consume one entire processor if the process is CPU bound. If the system also runs FSS, regular user processes compete for the three remaining CPUs that are not being used by the RT process. Note that the RT process might not use the CPU continuously. When the RT process is idle, FSS utilizes all four processors.

You can type the following command to find out which scheduling classes the processor sets are running in and ensure that each processor set is configured to run either TS, IA, FX, or FSS processes.


```
$ ps -ef -o pset,class | grep -v CLS | sort | uniq
1 FSS
1 SYS
2 TS
2 RT
3 FX
```

Setting the Scheduling Class for the System

To set the default scheduling class for the system, see [“How to Set the Scheduler Class” on page 109](#) and `dispadm(1M)`. To move running processes into a different scheduling class, see [“Configuring the FSS” on page 109](#) and `priocntl(1)`.

Scheduling Class in a Zones Environment

If you are using a zones environment on your system, the non-global zones use the default scheduling class for the system. If the system is updated with a new default scheduling class setting, non-global zones obtain the new setting when booted or rebooted.

It is also possible to set the scheduling class for a non-global zone through the dynamic resource pools facility. An association made through a resource pool will override the system setting. See [Step 6 in “How to Configure the Zone” on page 226](#) for information about setting the default scheduling class for a non-global zone by using the pools facility.

For information about moving running processes into a different scheduling class without changing the default scheduling class and rebooting, see [Table 25-2](#) the `priocntl(1)` man page.

Commands Used With FSS

The commands that are shown in the following table provide the primary administrative interface to the fair share scheduler.

Command Reference	Description
<code>priocntl(1)</code>	Displays or sets scheduling parameters of specified processes, moves running processes into a different scheduling class.
<code>ps(1)</code>	Lists information about running processes, identifies in which scheduling classes processor sets are running.
<code>dispadm(1M)</code>	Sets the default scheduler for the system. Also used to examine and tune the FSS scheduler's time quantum value.
<code>FSS(7)</code>	Describes the fair share scheduler (FSS).

Administering the Fair Share Scheduler (Tasks)

This chapter describes how to use the fair share scheduler (FSS).

For an overview of the FSS, see [Chapter 8](#).

Administering the Fair Share Scheduler (Task Map)

Task	Description	For Information
Monitor CPU usage.	Monitor the CPU usage of projects, and projects in processor sets.	“Monitoring the FSS” on page 108
Set the default scheduler class.	Make a scheduler such as the FSS the default scheduler for the system.	“How to Set the Scheduler Class” on page 109
Move running processes from one scheduler class to a different scheduling class, such as the FSS class.	Manually move processes from one scheduling class to another scheduling class without changing the default scheduling class and rebooting.	“How to Manually Move Processes From the TS Class Into the FSS Class” on page 109

Task	Description	For Information
Move all running processes from all scheduling classes to a different scheduling class, such as the FSS class.	Manually move processes in all scheduling classes to another scheduling class without changing the default scheduling class and rebooting.	“How to Manually Move Processes From All User Classes Into the FSS Class” on page 110
Move a project’s processes into a different scheduling class, such as the FSS class.	Manually move a project’s processes from their current scheduling class to a different scheduling class.	“How to Manually Move a Project’s Processes Into the FSS Class” on page 110
Examine and tune FSS parameters.	Tune the scheduler’s time quantum value. <i>Time quantum</i> is the amount of time that a thread is allowed to run before it must relinquish the processor.	“How to Tune Scheduler Parameters” on page 110

Monitoring the FSS

You can use the `prstat` command described in the `prstat(1M)` man page to monitor CPU usage by active projects.

You can use the extended accounting data for tasks to obtain per-project statistics on the amount of CPU resources that are consumed over longer periods. See [Chapter 4](#) for more information.

▼ How to Monitor System CPU Usage by Projects

- To monitor the CPU usage of projects that run on the system, use the `prstat` command with the `-J` option.

```
% prstat -J
```

▼ How to Monitor CPU Usage by Projects in Processor Sets

- To monitor the CPU usage of projects on a list of processor sets, type:

```
% prstat -J -C pset-list
```

where *pset-list* is a list of processor set IDs that are separated by commas.

Configuring the FSS

The same commands that you use with other scheduling classes in the Solaris system can be used with FSS. You can set the scheduler class, configure the scheduler's tunable parameters, and configure the properties of individual processes.

▼ How to Set the Scheduler Class

Use the `dispadmin` command to set FSS as the default scheduler for the system.

- 1. Become superuser or assume an equivalent role.**

Roles contain authorizations and privileged commands. For more information about roles, see “Using the Solaris Management Tools With RBAC (Task Map)” in *System Administration Guide: Basic Administration*.

- 2. Set the default scheduler for the system, for example, FSS.**

```
# dispadmin -d FSS
```

This change takes effect on the next reboot. After reboot, every process on the system runs in the FSS scheduling class.

▼ How to Manually Move Processes From the TS Class Into the FSS Class

You can manually move processes from one scheduling class to another scheduling class without changing the default scheduling class and rebooting. This procedure shows how to manually move processes from the TS scheduling class into the FSS scheduling class.

- 1. Become superuser or assume an equivalent role.**

Roles contain authorizations and privileged commands. For more information about roles, see “Using the Solaris Management Tools With RBAC (Task Map)” in *System Administration Guide: Basic Administration*.

- 2. Move the `init` process (pid 1) into the FSS scheduling class.**

```
# priocntl -s -c FSS -i pid 1
```

- 3. Move all processes from the TS scheduling class into the FSS scheduling class.**

```
# priocntl -s -c FSS -i class TS
```

All processes again run in the TS scheduling class after reboot.

▼ How to Manually Move Processes From All User Classes Into the FSS Class

You might be using a default class other than TS. For example, your system might be running a window environment that uses the IA class by default. You can manually move all processes into the FSS scheduling class without changing the default scheduling class and rebooting.

1. **Become superuser or assume an equivalent role.**

Roles contain authorizations and privileged commands. For more information about roles, see “Using the Solaris Management Tools With RBAC (Task Map)” in *System Administration Guide: Basic Administration*.

2. **Move the `init` process (pid 1) into the FSS scheduling class.**

```
# priocntl -s -c FSS -i pid 1
```

3. **Move all processes from their current scheduling classes into the FSS scheduling class.**

```
# priocntl -s -c FSS -i all
```

All processes again run in the default scheduling class after reboot.

▼ How to Manually Move a Project’s Processes Into the FSS Class

You can manually move a project’s processes from their current scheduling class to the FSS scheduling class.

1. **Become superuser or assume an equivalent role.**

Roles contain authorizations and privileged commands. For more information about roles, see “Using the Solaris Management Tools With RBAC (Task Map)” in *System Administration Guide: Basic Administration*.

2. **Move processes that run in project ID 10 to the FSS scheduling class.**

```
# priocntl -s -c FSS -i projid 10
```

The project’s processes again run in the default scheduling class after reboot.

How to Tune Scheduler Parameters

You can use the `dispadm` command to examine and tune the FSS scheduler’s time quantum value. *Time quantum* is the amount of time that a thread is allowed to run before it must relinquish the processor.

To display the current time quantum for the FSS scheduler, type:

```
$ dispadmin -c FSS -g
#
# Fair Share Scheduler Configuration
#
RES=1000
#
# Time Quantum
#
QUANTUM=110
```

When you use the `-g` option, you can also use the `-r` option to specify the resolution that is used for printing time quantum values. If no resolution is specified, time quantum values are displayed in milliseconds by default. Type the following:

```
$ dispadmin -c FSS -g -r 100
#
# Fair Share Scheduler Configuration
#
RES=100
#
# Time Quantum
#
QUANTUM=11
```

To set scheduling parameters for the FSS scheduling class, use `dispadmin -s`. The values in *file* must be in the format output by the `-g` option. These values overwrite the current values in the kernel. Type the following:

```
$ dispadmin -c FSS -s file
```


Physical Memory Control Using the Resource Capping Daemon (Overview)

The resource capping daemon `rcapd` regulates physical memory consumption by processes running in projects that have resource caps defined.

The following topics are covered in this chapter.

- “Introduction to the Resource Capping Daemon” on page 113
- “How Resource Capping Works” on page 114
- “Attribute to Limit Physical Memory Usage” on page 115
- “`rcapd` Configuration” on page 115
- “Monitoring Resource Utilization With `rcapstat`” on page 119
- “Commands Used With `rcapd`” on page 121

For procedures using the `rcapd` feature, see [Chapter 11](#).

Introduction to the Resource Capping Daemon

A resource *cap* is an upper bound placed on the consumption of a resource, such as physical memory. Per-project physical memory caps are supported.

The resource capping daemon and its associated utilities provide mechanisms for physical memory resource cap enforcement and administration.

Like the resource control, the resource cap can be defined by using attributes of project entries in the `project` database. However, while resource controls are synchronously enforced by the kernel, resource caps are asynchronously enforced at the user level by the resource capping daemon. With asynchronous enforcement, a small delay occurs as a result of the sampling interval used by the daemon.

For information about `rcapd`, see the `rcapd(1M)` man page. For information about projects and the `project` database, see [Chapter 2](#) and the `project(4)` man page. For information about resource controls, see [Chapter 6](#).

Note – If you are using `rcapd` in a zones environment, you must add a `project` entry and configure the daemon in each zone where you want the daemon to run. `rcapd` will not act on processes in zones other than the one in which it is running.

How Resource Capping Works

The daemon repeatedly samples the resource utilization of projects that have physical memory caps. The sampling interval used by the daemon is specified by the administrator. See [“Determining Sample Intervals” on page 119](#) for additional information. When the system’s physical memory utilization exceeds the threshold for cap enforcement, and other conditions are met, the daemon takes action to reduce the resource consumption of projects with memory caps to levels at or below the caps.

The virtual memory system divides physical memory into segments known as pages. Pages are the fundamental unit of physical memory in the Solaris memory management subsystem. To read data from a file into memory, the virtual memory system reads in one page at a time, or *pages in* a file. To reduce resource consumption, the daemon can *page out*, or relocate, infrequently used pages to a swap device, which is an area outside of physical memory.

The daemon manages physical memory by regulating the size of a project workload’s resident set relative to the size of its working set. The resident set is the set of pages that are resident in physical memory. The working set is the set of pages that the workload actively uses during its processing cycle. The working set changes over time, depending on the process’s mode of operation and the type of data being processed. Ideally, every workload has access to enough physical memory to enable its working set to remain resident. However, the working set can also include the use of secondary disk storage to hold the memory that does not fit in physical memory.

Only one instance of `rcapd` can run at any given time.

Attribute to Limit Physical Memory Usage

To define a physical memory resource cap for a project, establish a resident set size (RSS) cap by adding this attribute to the `project` database entry:

`rcap.max-rss` The total amount of physical memory, in bytes, that is available to processes in the project.

For example, the following line in the `/etc/project` file sets an RSS cap of 10 gigabytes for a project named `db`.

```
db:100::db,root::rcap.max-rss=10737418240
```

Note – The system might round the specified cap value to a page size.

You can use the `projmod` command to set the `rcap.max-rss` attribute in the `/etc/project` file:

```
# projmod -s -K rcap.max-rss=10GB db
```

The `/etc/project` file then contains the line:

```
db:100::db,root::rcap.max-rss=10737418240
```

rcapd Configuration

You use the `rcapadm` command to configure the resource capping daemon. You can perform the following actions:

- Set the threshold value for cap enforcement
- Set intervals for the operations performed by `rcapd`
- Enable or disable resource capping
- Display the current status of the configured resource capping daemon

To configure the daemon, you must have superuser privileges or have the Process Management profile in your list of profiles. The Process Management role and the System Administrator role both include the Process Management profile.

Configuration changes can be incorporated into `rcapd` according to the configuration interval (see [“`rcapd` Operation Intervals” on page 118](#)) or on demand by sending a `SIGHUP` (see the `kill(1)` man page).

If used without arguments, `rcapadm` displays the current status of the resource capping daemon if it has been configured.

The following subsections discuss cap enforcement, cap values, and `rcapd` operation intervals.

Memory Cap Enforcement Threshold

The *memory cap enforcement threshold* is the percentage of physical memory utilization on the system that triggers cap enforcement. When the system exceeds this utilization, caps are enforced. The physical memory used by applications and the kernel is included in this percentage. The percentage of utilization determines the way in which memory caps are enforced.

To enforce caps, memory can be paged out from project workloads.

- Memory can be paged out to reduce the size of the portion of memory that is over its cap for a given workload.
- Memory can be paged out to reduce the proportion of physical memory used that is over the memory cap enforcement threshold on the system.

A workload is permitted to use physical memory up to its cap. A workload can use additional memory as long as the system’s memory utilization stays below the memory cap enforcement threshold.

To set the value for cap enforcement, see [“How to Set the Memory Cap Enforcement Threshold” on page 124](#).

Determining Cap Values

If a project cap is set too low, there might not be enough memory for the workload to proceed effectively under normal conditions. The paging that occurs because the workload requires more memory has a negative effect on system performance.

Projects that have caps set too high can consume available physical memory before their caps are exceeded. In this case, physical memory is effectively managed by the kernel and not by `rcapd`.

In determining caps on projects, consider these factors.

Impact on I/O system	<p>The daemon can attempt to reduce a project workload's physical memory usage whenever the sampled usage exceeds the project's cap. During cap enforcement, the swap devices and other devices that contain files that the workload has mapped are used. The performance of the swap devices is a critical factor in determining the performance of a workload that routinely exceeds its cap. The execution of the workload is similar to running it on a machine with the same amount of physical memory as the workload's cap.</p>
Impact on CPU usage	<p>The daemon's CPU usage varies with the number of processes in the project workloads it is capping and the sizes of the workloads' address spaces.</p> <p>A small portion of the daemon's CPU time is spent sampling the usage of each workload. Adding processes to workloads increases the time spent sampling usage.</p> <p>Another portion of the daemon's CPU time is spent enforcing caps when they are exceeded. The time spent is proportional to the amount of virtual memory involved. CPU time spent increases or decreases in response to corresponding changes in the total size of a workload's address space. This information is reported in the <code>vm</code> column of <code>rcapstat</code> output. See "Monitoring Resource Utilization With <code>rcapstat</code>" on page 119 and the <code>rcapstat(1)</code> man page for more information.</p>
Reporting on shared memory	<p>The daemon cannot determine which pages of memory are shared with other processes or which are mapped multiple times within the same process. Since <code>rcapd</code> assumes that each page is unique, this results in a discrepancy between the reported (estimated) RSS and the actual RSS.</p> <p>Certain workloads, such as databases, use shared memory extensively. For these workloads, you can sample a project's regular usage to determine a suitable initial cap value. Use output from the <code>prstat</code> command with the <code>-J</code> option. See the <code>prstat(1M)</code> man page.</p>

rcapd Operation Intervals

You can tune the intervals for the periodic operations performed by `rcapd`.

All intervals are specified in seconds. The `rcapd` operations and their default interval values are described in the following table.

Operation	Default Interval Value in Seconds	Description
<code>scan</code>	15	Number of seconds between scans for processes that have joined or left a project workload. Minimum value is 1 second.
<code>sample</code>	5	Number of seconds between samplings of resident set size and subsequent cap enforcements. Minimum value is 1 second.
<code>report</code>	5	Number of seconds between updates to paging statistics. If set to 0, statistics are not updated, and output from <code>rcapstat</code> is not current.
<code>config</code>	60	Number of seconds between reconfigurations. In a reconfiguration event, <code>rcapadm</code> reads the configuration file for updates, and scans the project database for new or revised project caps. Sending a <code>SIGHUP</code> to <code>rcapd</code> causes an immediate reconfiguration.

To tune intervals, see [“How to Set Operation Intervals”](#) on page 125.

Determining `rcapd` Scan Intervals

The scan interval controls how often `rcapd` looks for new processes. On systems with many processes running, the scan through the list takes more time, so it might be preferable to lengthen the interval in order to reduce the overall CPU time spent. However, the scan interval also represents the minimum amount of time that a process must exist to be attributed to a capped workload. If there are workloads that run many short-lived processes, `rcapd` might not attribute the processes to a workload if the scan interval is lengthened.

Determining Sample Intervals

The sample interval configured with `rcapadm` is the shortest amount of time `rcapd` waits between sampling a workload's usage and enforcing the cap if it is exceeded. If you reduce this interval, `rcapd` will, under most conditions, enforce caps more frequently, possibly resulting in increased I/O due to paging. However, a shorter sample interval can also lessen the impact that a sudden increase in a particular workload's physical memory usage might have on other workloads. The window between samplings, in which the workload can consume memory unhindered and possibly take memory from other capped workloads, is narrowed.

If the sample interval specified to `rcapstat` is shorter than the interval specified to `rcapd` with `rcapadm`, the output for some intervals can be zero. This situation occurs because `rcapd` does not update statistics more frequently than the interval specified with `rcapadm`. The interval specified with `rcapadm` is independent of the sampling interval used by `rcapstat`.

Monitoring Resource Utilization With `rcapstat`

Use `rcapstat` to monitor the resource utilization of capped projects. To view an example `rcapstat` report, see [“Producing Reports With `rcapstat`”](#) on page 126.

You can set the sampling interval for the report and specify the number of times that statistics are repeated.

interval Specifies the sampling interval in seconds. The default interval is 5 seconds.

count Specifies the number of times that the statistics are repeated. By default, `rcapstat` reports statistics until a termination signal is received or until the `rcapd` process exits.

The paging statistics in the first report issued by `rcapstat` show the activity since the daemon was started. Subsequent reports reflect the activity since the last report was issued.

The following table defines the column headings in an `rcapstat` report.

<code>rcapstat</code> Column Headings	Description
<code>id</code>	The project ID of the capped project.

<code>rcapstat</code> Column Headings	Description
<code>project</code>	The project name.
<code>nproc</code>	The number of processes in the project.
<code>vm</code>	The total amount of virtual memory size used by processes in the project, in kilobytes (K), megabytes (M), or gigabytes (G).
<code>rss</code>	The estimated amount of the total resident set size (RSS) of the processes in the project, in kilobytes (K), megabytes (M), or gigabytes (G), not accounting for pages that are shared.
<code>cap</code>	The RSS cap defined for the project. See “Attribute to Limit Physical Memory Usage” on page 115 or the <code>rcapd(1M)</code> man page for information about how to specify memory caps.
<code>at</code>	The total amount of memory that <code>rcapd</code> attempted to page out since the last <code>rcapstat</code> sample.
<code>avgat</code>	The average amount of memory that <code>rcapd</code> attempted to page out during each sample cycle that occurred since the last <code>rcapstat</code> sample. The rate at which <code>rcapd</code> samples collection RSS can be set with <code>rcapadm</code> . See “rcapd Operation Intervals” on page 118 .
<code>pg</code>	The total amount of memory that <code>rcapd</code> successfully paged out since the last <code>rcapstat</code> sample.
<code>avgpg</code>	An estimate of the average amount of memory that <code>rcapd</code> successfully paged out during each sample cycle that occurred since the last <code>rcapstat</code> sample. The rate at which <code>rcapd</code> samples process RSS sizes can be set with <code>rcapadm</code> . See “rcapd Operation Intervals” on page 118 .

Commands Used With `rcapd`

Command Reference	Description
<code>rcapstat(1)</code>	Monitors the resource utilization of capped projects.
<code>rcapadm(1M)</code>	Configures the resource capping daemon, displays the current status of the resource capping daemon if it has been configured, and enables or disables resource capping
<code>rcapd(1M)</code>	The resource capping daemon.

Administering the Resource Capping Daemon (Tasks)

This chapter contains procedures for configuring and using the resource capping daemon `rcapd`.

For an overview of `rcapd`, see [Chapter 10](#).

Configuring and Using the Resource Capping Daemon (Task Map)

Task	Description	For Instructions
Set the memory cap enforcement threshold.	Configure a cap that will be enforced when the physical memory available to processes is low.	“How to Set the Memory Cap Enforcement Threshold” on page 124
Set the operation interval.	The interval is applied to the periodic operations performed by the resource capping daemon.	“How to Set Operation Intervals” on page 125
Enable resource capping.	Activate resource capping on your system.	“How to Enable Resource Capping” on page 125
Disable resource capping.	Deactivate resource capping on your system.	“How to Disable Resource Capping” on page 126
Report cap and project information.	View example commands for producing reports.	“Reporting Cap and Project Information” on page 126

Task	Description	For Instructions
Monitor a project's resident set size.	Produce a report on the resident set size of a project.	"Monitoring the RSS of a Project" on page 127
Determine a project's working set size.	Produce a report on the working set size of a project.	"Determining the Working Set Size of a Project" on page 128
Report on memory utilization and memory caps.	Print a memory utilization and cap enforcement line at the end of the report for each interval.	"Reporting Memory Utilization and the Memory Cap Enforcement Threshold" on page 130

Administering the Resource Capping Daemon With `rcapadm`

This section contains procedures for configuring the resource capping daemon with `rcapadm`. See "[rcapd Configuration](#)" on page 115 and the `rcapadm(1M)` man page for more information.

If used without arguments, `rcapadm` displays the current status of the resource capping daemon if it has been configured.

▼ How to Set the Memory Cap Enforcement Threshold

Caps can be configured so that they will not be enforced until the physical memory available to processes is low. See "[Memory Cap Enforcement Threshold](#)" on page 116 for more information.

The minimum (and default) value is 0, which means that memory caps are always enforced. To set a different minimum, follow this procedure.

1. Become superuser, or assume a role that includes the Process Management profile.

The System Administrator role includes the Process Management profile. For information on how to create the role and assign the role to a user, see "Managing RBAC (Task Map)" in *System Administration Guide: Security Services*.

2. Use the `-c` option of `rcapadm` to set a different physical memory utilization value for memory cap enforcement.

```
# rcapadm -c percent
```

percent is in the range 0 to 100. Higher values are less restrictive. A higher value means capped project workloads can execute without having caps enforced until the system's memory utilization exceeds this threshold.

To display the current physical memory utilization and the cap enforcement threshold, see [“Reporting Memory Utilization and the Memory Cap Enforcement Threshold”](#) on page 130.

▼ How to Set Operation Intervals

[“rcapd Operation Intervals”](#) on page 118 contains information about the intervals for the periodic operations performed by `rcapd`. To set operation intervals using `rcapadm`, follow this procedure.

1. **Become superuser, or assume a role that includes the Process Management profile.**

The System Administrator role includes the Process Management profile. For information on how to create the role and assign the role to a user, see [“Managing RBAC \(Task Map\)”](#) in *System Administration Guide: Security Services*.

2. **Use the `-i` option to set interval values.**

```
# rcapadm -i interval=value, ..., interval=value
```

All interval values are specified in seconds.

▼ How to Enable Resource Capping

There are two ways to enable resource capping on your system.

1. **Become superuser, or assume a role that includes the Process Management profile.**

The System Administrator role includes the Process Management profile. For information on how to create the role and assign the role to a user, see [“Managing RBAC \(Task Map\)”](#) in *System Administration Guide: Security Services*.

2. **Enable the resource capping daemon in one of the following ways:**

- To enable the resource capping daemon so that it will be started now and also be started each time the system is booted, type:

```
# rcapadm -E
```

- To enable the resource capping daemon at boot without starting it now, also specify the `-n` option:

```
# rcapadm -n -E
```

▼ How to Disable Resource Capping

There are two ways to disable resource capping on your system.

1. **Become superuser, or assume a role that includes the Process Management profile.**

The System Administrator role includes the Process Management profile. For information on how to create the role and assign the role to a user, see “Managing RBAC (Task Map)” in *System Administration Guide: Security Services*.

2. **Disable the resource capping daemon in one of the following ways:**

- To disable the resource capping daemon so that it will be stopped now and not be started when the system is booted, type:

```
# rcapadm -D
```

- To disable the resource capping daemon without stopping it, also specify the `-n` option:

```
# rcapadm -n -D
```

Note – Use `rcapadm -D` to safely disable `rcapd`. If the daemon is killed (see the `kill(1)` man page), processes might be left in a stopped state and need to be manually restarted. To resume a process running, use the `prun` command. See the `prun(1)` man page for more information.

Producing Reports With `rcapstat`

Use `rcapstat` to report resource capping statistics. “[Monitoring Resource Utilization With `rcapstat`](#)” on page 119 explains how to use the `rcapstat` command to generate reports. That section also describes the column headings in the report. The `rcapstat(1)` man page also contains this information.

The following subsections use examples to illustrate how to produce reports for specific purposes.

Reporting Cap and Project Information

In this example, caps are defined for two projects associated with two users. `user1` has a cap of 50 megabytes, and `user2` has a cap of 10 megabytes.

The following command produces five reports at 5-second sampling intervals.

```

user1machine% rcapstat 5 5
  id project nproc  vm  rss  cap  at avgat  pg avgpg
112270 user1  24  123M  35M  50M  50M  0K 3312K  0K
 78194 user2  1  2368K  1856K  10M  0K  0K  0K  0K
  id project nproc  vm  rss  cap  at avgat  pg avgpg
112270 user1  24  123M  35M  50M  0K  0K  0K  0K
 78194 user2  1  2368K  1856K  10M  0K  0K  0K  0K
  id project nproc  vm  rss  cap  at avgat  pg avgpg
112270 user1  24  123M  35M  50M  0K  0K  0K  0K
 78194 user2  1  2368K  1928K  10M  0K  0K  0K  0K
  id project nproc  vm  rss  cap  at avgat  pg avgpg
112270 user1  24  123M  35M  50M  0K  0K  0K  0K
 78194 user2  1  2368K  1928K  10M  0K  0K  0K  0K
  id project nproc  vm  rss  cap  at avgat  pg avgpg
112270 user1  24  123M  35M  50M  0K  0K  0K  0K
 78194 user2  1  2368K  1928K  10M  0K  0K  0K  0K

```

The first three lines of output constitute the first report, which contains the cap and project information for the two projects and paging statistics since rcapd was started. The at and pg columns are a number greater than zero for user1 and zero for user2, which indicates that at some time in the daemon's history, user1 exceeded its cap but user2 did not.

The subsequent reports show no significant activity.

Monitoring the RSS of a Project

The following example shows project user1, which has an RSS in excess of its RSS cap.

The following command produces five reports at 5-second sampling intervals.

```

user1machine% rcapstat 5 5
  id project nproc  vm  rss  cap  at avgat  pg avgpg
376565 user1  3  6249M  6144M  6144M  690M  220M  5528K  2764K
376565 user1  3  6249M  6144M  6144M  0M  131M  4912K  1637K
376565 user1  3  6249M  6171M  6144M  27M  147M  6048K  2016K
376565 user1  3  6249M  6146M  6144M  4872M  174M  4368K  1456K
376565 user1  3  6249M  6156M  6144M  12M  161M  3376K  1125K

```

The user1 project has three processes that are actively using physical memory. The positive values in the pg column indicate that rcapd is consistently paging out memory as it attempts to meet the cap by lowering the physical memory utilization of the project's processes. However, rcapd does not succeed in keeping the RSS below the cap value. This is indicated by the varying rss values that do not show a corresponding decrease. As soon as memory is paged out, the workload uses it again and the RSS count goes back up. This means that all of the project's resident memory is being actively used and the working set size (WSS) is greater than the cap. Thus, rcapd is forced to page out some of the working set to meet the cap. Under this condition, the system will continue to experience high page fault rates, and associated I/O, until one of the following occurs:

- The WSS becomes smaller.
- The cap is raised.
- The application changes its memory access pattern.

In this situation, shortening the sample interval might reduce the discrepancy between the RSS value and the cap value by causing `rcapd` to sample the workload and enforce caps more frequently.

Note – A page fault occurs when either a new page must be created or the system must copy in a page from a swap device.

Determining the Working Set Size of a Project

The following example is a continuation of the previous example, and it uses the same project.

The previous example shows that the `user1` project is using more physical memory than its cap allows. This example shows how much memory the project workload requires.

```
user1machine% rcapstat 5 5
  id project  nproc   vm   rss   cap   at avgat   pg  avgpg
376565  user1     3 6249M 6144M 6144M 690M   0K   689M   0K
376565  user1     3 6249M 6144M 6144M   0K   0K    0K   0K
376565  user1     3 6249M 6171M 6144M  27M   0K    27M   0K
376565  user1     3 6249M 6146M 6144M 4872K   0K  4816K   0K
376565  user1     3 6249M 6156M 6144M  12M   0K   12M   0K
376565  user1     3 6249M 6150M 6144M 5848K   0K  5816K   0K
376565  user1     3 6249M 6155M 6144M  11M   0K   11M   0K
376565  user1     3 6249M 6150M  10G   32K   0K   32K   0K
376565  user1     3 6249M 6214M  10G   0K   0K   0K   0K
376565  user1     3 6249M 6247M  10G   0K   0K   0K   0K
376565  user1     3 6249M 6247M  10G   0K   0K   0K   0K
376565  user1     3 6249M 6247M  10G   0K   0K   0K   0K
376565  user1     3 6249M 6247M  10G   0K   0K   0K   0K
376565  user1     3 6249M 6247M  10G   0K   0K   0K   0K
376565  user1     3 6249M 6247M  10G   0K   0K   0K   0K
```

Halfway through the cycle, the cap on the `user1` project was increased from 6 gigabytes to 10 gigabytes. This increase stops cap enforcement and allows the resident set size to grow, limited only by other processes and the amount of memory in the machine. The `rss` column might stabilize to reflect the project working set size (WSS), 6247M in this example. This is the minimum cap value that allows the project's processes to operate without continually incurring page faults.

The following two figures graphically show the effect `rcapd` has on `user1` while the cap is 6 gigabytes and 10 gigabytes. Every 5 seconds, corresponding to the sample interval, the RSS decreases and I/O increases as `rcapd` pages out some of the workload's memory. Shortly after the page out completes, the workload, needing

those pages, pages them back in as it continues running. This cycle repeats until the cap is raised to 10 gigabytes approximately halfway through the example, and the RSS stabilizes at 6.1 gigabytes. Since the workload's RSS is now below the cap, no more paging occurs. The I/O associated with paging stops as well, as the `vmstat` (see `vmstat(1M)`) or `iostat` (see `iostat(1M)`) commands would show. Thus, you can infer that the project required 6.1 gigabytes to perform the work it was doing at the time it was being observed.

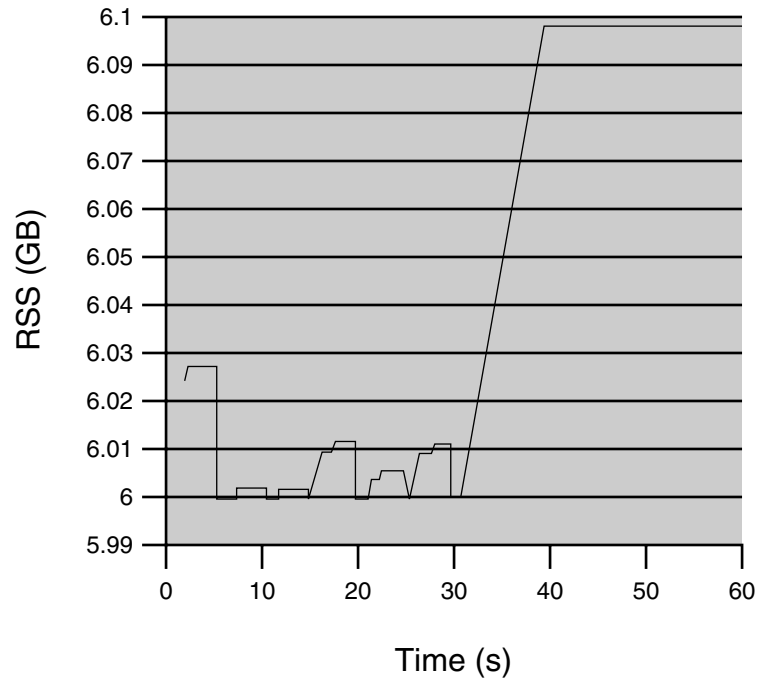
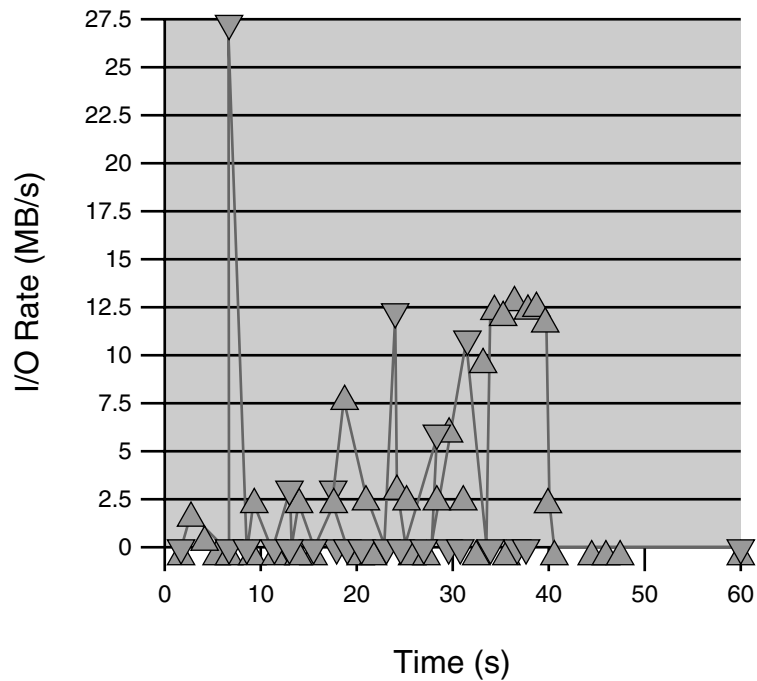


FIGURE 11-1 Stabilizing RSS Values After Raising the Cap of `user1` Higher Than `user1`'s WSS



▼ Page Outs
 ▲ Page Ins

FIGURE 11-2 Relationship Between Page Ins and Page Outs, and the Stabilization of I/O After user1's Cap Is Raised

Reporting Memory Utilization and the Memory Cap Enforcement Threshold

You can use the `-g` option of `rcapstat` to report the following:

- Current physical memory utilization as a percentage of physical memory installed on the system
- System memory cap enforcement threshold set by `rcapadm`

The `-g` option causes a memory utilization and cap enforcement line to be printed at the end of the report for each interval.

```
# rcapstat -g
id project nproc vm rss cap at avgat pg avgpg
376565 rcap 0 0K 0K 10G 0K 0K 0K 0K
```

```
physical memory utilization: 55%   cap enforcement threshold: 0%
  id project  nproc   vm  rss  cap   at avgat  pg  avgpg
376565  rcap      0   0K  0K  10G  0K   0K   0K   0K
physical memory utilization: 55%   cap enforcement threshold: 0%
```


Dynamic Resource Pools (Overview)

This chapter discusses resource pools, which are used for partitioning machine resources. Dynamic resource pools (DRPs) dynamically adjust each resource pool's resource allocation to meet established system goals.

The following topics are covered in this chapter:

- "Introduction to Resource Pools" on page 133
- "project.pool Attribute" on page 135
- "Resource Pools Used in Zones" on page 135
- "When to Use Pools" on page 136
- "Resource Pools Framework" on page 137
- "Implementing Pools on a System" on page 139
- "SPARC: Dynamic Reconfiguration Operations and Resource Pools" on page 140
- "Creating Pools Configurations" on page 140
- "Directly Manipulating the Dynamic Configuration" on page 141
- "poold Overview" on page 141
- "Configuration Constraints and Objectives" on page 142
- "poold Features That Can Be Configured" on page 147
- "How Dynamic Resource Allocation Works" on page 150
- "Commands Used With the Resource Pools Facility" on page 155

For procedures using this functionality, see Chapter 13.

Introduction to Resource Pools

Resource pools enable you to separate workloads so that workload consumption of certain resources does not overlap. This resource reservation helps to achieve predictable performance on systems with mixed workloads.

Resource pools provide a persistent configuration mechanism for processor set (pset) configuration and, optionally, scheduling class assignment.

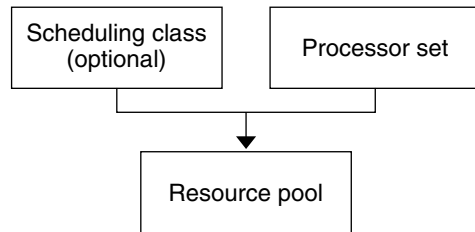


FIGURE 12-1 Resource Pool Framework

A pool can be thought of as a specific binding of the various resource sets that are available on your system. You can create pools that represent different kinds of possible resource combinations:

```

pool1: pset_default
pool2: pset1
pool3: pset1, pool.scheduler="FSS"
  
```

By grouping multiple partitions, pools provide a handle to associate with labeled workloads. Each project entry in the `/etc/project` file can have a single pool associated with that entry, which is specified using the `project.pool` attribute.

Resource pools provide a mechanism for dynamically adjusting each pool's resource allocation in response to system events and application load changes. DRPs simplify and reduce the number of decisions required from an administrator. Adjustments are automatically made to preserve the system performance goals specified by an administrator. The changes made to the configuration are logged. These features are primarily enacted through the resource controller `poold`, a system daemon that should always be active when dynamic resource allocation is required. Periodically, `poold` examines the load on the system and determines whether intervention is required to enable the system to maintain optimal performance with respect to resource consumption. The `poold` configuration is held in the `libpool` configuration. For more information on `poold`, see the `poold(1M)` man page.

When pools are enabled, a *default pool* and a *default processor set* form the base configuration. Additional user-defined pools and processor sets can be created and added to the configuration. A CPU can only belong to one processor set. User-defined pools and processor sets can be destroyed. The default pool and the default processor set cannot be destroyed.

The default pool has the `pool.default` property set to `true`. The default processor set has the `pset.default` property set to `true`. Thus, both the default pool and the default processor set can be identified even if their names have been changed.

The user-defined pools mechanism is primarily for use on large machines of more than four CPUs. However, small machines can still benefit from this functionality. On small machines, you can create pools that share noncritical resource partitions. The pools are separated only on the basis of critical resources.

project.pool Attribute

The `project.pool` attribute can be added to a project entry in the `/etc/project` file to associate a single pool with that entry. New work that is started on a project is bound to the appropriate pool. See [Chapter 2](#) for more information.

for example, you can use the `projmod` command to set the `project.pool` attribute for the project `sales` in the `/etc/project` file:

```
# projmod -s -K project.pool=mypool sales
```

Resource Pools Used in Zones

On a system that has zones enabled, a non-global zone can be associated with one resource pool. Moreover, you cannot bind individual processes in non-global zones to a different pool by using the `poolbind` command from the global zone. To associate a non-global zone with a pool, see [“Configuring, Verifying, and Committing a Zone” on page 226](#). Note that if you set a scheduling class for a pool and you associate a non-global zone with that pool, the zone uses that scheduling class by default.

The scope of an executing instance of `poold` is limited to the global zone.

The `poolstat` utility run in a non-global zone displays only information about the pool associated with the zone. The `pooladm` command run without arguments in a non-global zone displays only information about the pool associated with the zone.

For information about resource pool commands, see [“Commands Used With the Resource Pools Facility” on page 155](#).

When to Use Pools

Resource pools offer a versatile mechanism that can be applied to many administrative scenarios.

Batch compute server	Use pools functionality to split a server into two pools. One pool is used for login sessions and interactive work by timesharing users. The other pool is used for jobs that are submitted through the batch system.
Application or database server	Partition the resources for interactive applications in accordance with the applications' requirements.
Turning on applications in phases	<p>Set user expectations.</p> <p>You might initially deploy a machine that is running only a fraction of the services that the machine is ultimately expected to deliver. User difficulties can occur if reservation-based resource management mechanisms are not established when the machine comes online.</p> <p>For example, the fair share scheduler optimizes CPU utilization. The response times for a machine that is running only one application can be misleadingly fast. Users will not see these response times with multiple applications loaded. By using separate pools for each application, you can place a ceiling on the number of CPUs available to each application before you deploy all applications.</p>
Complex timesharing server	<p>Partition a server that supports large user populations. Server partitioning provides an isolation mechanism that leads to a more predictable per-user response.</p> <p>By dividing users into groups that bind to separate pools, and using the fair share scheduling (FSS) facility, you can tune CPU allocations to favor sets of users that have priority. This assignment can be based on user role, accounting chargeback, and so forth.</p>

Workloads that change seasonally	Use resource pools to adjust to changing demand. Your site might experience predictable shifts in workload demand over long periods of time, such as monthly, quarterly, or annual cycles. If your site experiences these shifts, you can alternate between multiple pools configurations by invoking <code>pooladm</code> from a <code>cron</code> job. (See “Resource Pools Framework” on page 137.)
Real-time applications	Create a real-time pool by using the RT scheduler and designated processor resources.
System utilization	Enforce system goals that you establish. Use the automated pools daemon feature to identify available resources and then monitor workloads to detect when your specified objectives are no longer being satisfied. The daemon can take corrective action if possible, or the condition can be logged.

Resource Pools Framework

The `/etc/pooladm.conf` configuration file describes the static pools configuration. A static configuration represents the way in which an administrator would like a system to be configured with respect to resource pools functionality. An alternate file name can be specified.

The kernel holds information about the disposition of resources within the resource pools framework. This is known as the dynamic configuration, and it represents the resource pools functionality for a particular system at a point in time. The dynamic configuration can be viewed by using the `pooladm` command. Note that the order in which properties are displayed for pools and resource sets can vary. Modifications to the dynamic configuration are made indirectly, by applying a static configuration file, or directly, by using the `poolcfg` command with the `-d` option.

More than one static pools configuration file can exist, for activation at different times. You can alternate between multiple pools configurations by invoking `pooladm` from a `cron` job. (See the `cron(1M)` man page.)

By default, the resource pools framework is not active. Resource pools must be enabled to create or modify the dynamic configuration. Static configuration files can be manipulated with the `poolcfg` or `libpool` commands even if the resource pools framework is disabled. Static configuration files cannot be created if the pools facility is not active. For more information on the configuration file, see [“Creating Pools Configurations” on page 140](#).

The commands used with resource pools and the `poold` system daemon are described in the following man pages:

- `pooladm(1M)`
- `poolcfg(1M)`
- `poold(1M)`
- `libpool(3LIB)`

`/etc/pooladm.conf` Contents

The static configuration held at `/etc/pooladm.conf` is a special static configuration. When a system boots, if this file exists, then the resource pools framework is enabled and this static configuration is applied to the system.

All resource pool configurations, including the dynamic configuration, can contain the following elements.

<code>system</code>	Properties affecting the total behavior of the system
<code>pool</code>	A resource pool definition
<code>pset</code>	A processor set definition
<code>cpu</code>	A processor definition

All of these elements have properties that can be manipulated to alter the state and behavior of the resource pools framework. For example, the `pool` property `pool.importance` indicates the relative importance of a given pool. This property is used for possible resource dispute resolution. For more information, see `libpool(3LIB)`.

Pools Properties

The pools facility supports named, typed properties that can be placed on a pool, resource, or component. Administrators can store additional properties on the various pool elements. A property namespace similar to the `project` attribute is used.

For example, the following comment indicates that a given `pset` is associated with a particular `Datatree` database.

Datatree,pset.dbname=warehouse

For additional information about property types, see [“poold Properties”](#) on page 146.

Note – A number of special properties are reserved for internal use and cannot be set or removed. See the `libpool(3LIB)` man page for more information.

Implementing Pools on a System

User-defined pools can be implemented on a system by using one of these methods.

- When the Solaris software boots, an `init` script checks to see if the `/etc/pooladm.conf` file exists. If this file is found, then `pooladm` is invoked to make this configuration the active pools configuration. The system creates a dynamic configuration to reflect the organization that is requested in `/etc/pooladm.conf`, and the machine’s resources are partitioned accordingly.
- When the Solaris system is running, a pools configuration can either be activated if it is not already present, or modified by using the `pooladm` command. By default, the `pooladm` command operates on `/etc/pooladm.conf`. However, you can optionally specify an alternate location and file name, and use that file to update the pools configuration.

The `poold` resource controller is started with the pools facility.

For information about enabling and disabling resource pools, see [“Enabling and Disabling the Pools Facility”](#) on page 159. The pools facility cannot be disabled when there are user-defined pools or resources in use.

To configure resource pools, you must have superuser privileges or have the Process Management profile in your list of profiles. The Process Management role and the System Administrator role both include the Process Management profile.

SPARC: Dynamic Reconfiguration Operations and Resource Pools

Dynamic Reconfiguration (DR) enables you to reconfigure hardware while the system is running. A DR operation can increase, reduce, or have no effect on a given type of resource. Because DR can affect available resource amounts, the pools facility must be included in these operations. When a DR operation is initiated, the pools framework acts to validate the configuration.

If the DR operation can proceed without causing the current pools configuration to become invalid, then the private configuration file is updated. An invalid configuration is one that cannot be supported by the available resources.

If the DR operation would cause the pools configuration to be invalid, then the operation fails and you are notified by a message to the message log. If you want to force the configuration to completion, you must use the DR force option. The pools configuration is then modified to comply with the new resource configuration. For information on the DR process and the force option, see the dynamic reconfiguration user guide for your Sun hardware.

It is possible for a partition to move out of `poold` control while the daemon is active. For more information, see [“Identifying a Resource Shortage” on page 152](#).

Creating Pools Configurations

The configuration file contains a description of the pools to be created on the system. The file describes the elements that can be manipulated.

- system
- pool
- pset
- cpu

See `poolcfg(1M)` for more information on elements that be manipulated.

When pools are enabled, you can create a structured `/etc/pooladm.conf` file in two ways.

- You can use the `pooladm` command with the `-s` option to discover the resources on the current system and place the results in a configuration file.

This method is preferred. All active resources and components on the system that are capable of being manipulated by the pools facility are recorded. The resources include existing processor set configurations. You can then modify the configuration to rename the processor sets or to create additional pools if necessary.

- You can use the `poolcfg` command with the `-c` option and the `discover` or `create system name` subcommands to create a new pools configuration.

These options are maintained for backward compatibility with the previous release.

Use `poolcfg` or `libpool` to modify the `/etc/pooladm.conf` file. Do not directly edit this file.

Directly Manipulating the Dynamic Configuration

It is possible to directly manipulate CPU resource types in the dynamic configuration by using the `poolcfg` command with the `-d` option. There are two methods used to transfer resources.

- You can make a general request to transfer any available identified resources between sets.
- You can transfer resources with specific IDs to a target set. Note that the system IDs associated with resources can change when the resource configuration is altered or after a system reboot.

For an example, see [“Transferring Resources” on page 169](#).

Note that the resource transfer might trigger action from `poold`. See [“`poold` Overview” on page 141](#) for more information.

`poold` Overview

The pools resource controller, `poold`, uses system targets and observable statistics to preserve the system performance goals that you specify. This system daemon should always be active when dynamic resource allocation is required.

The `pool`d resource controller identifies available resources and then monitors workloads to determine when the system usage objectives are no longer being met. `pool`d then considers alternative configurations in terms of the objectives, and remedial action is taken. If possible, the resources are reconfigured so that objectives can be met. If this action is not possible, the daemon logs that user-specified objectives can no longer be achieved. Following a reconfiguration, the daemon resumes monitoring workload objectives.

`pool`d maintains a decision history that it can examine. The decision history is used to eliminate reconfigurations that historically did not show improvements.

Note that a reconfiguration can also be triggered asynchronously if the workload objectives are changed or if the resources available to the system are modified.

Stopping `pool`d

If, for any reason, dynamic resource allocation is not required, `pool`d can be stopped with the `SIGQUIT` or the `SIGTERM` signal. Either of these signals causes `pool`d to terminate gracefully.

Reconfiguring `pool`d

`pool`d will automatically detect changes in the resource or pools configuration. However, you can also force a reconfiguration to occur by using the `SIGHUP` signal.

Configuration Constraints and Objectives

When making changes to a configuration, `pool`d acts on directions that you provide. You specify these directions as a series of constraints and objectives. `pool`d uses your specifications to determine the relative value of different configuration possibilities in relation to the existing configuration. `pool`d then changes the resource assignments of the current configuration to generate new candidate configurations.

Configuration Constraints

Constraints affect the range of possible configurations by eliminating some of the potential changes that could be made to a configuration. The following constraints, which are specified in the `libpool` configuration, are available.

- The minimum and maximum CPU allocations
- Pinned components that are not available to be moved from a set

See the `libpool(3LIB)` man page and “[Pools Properties](#)” on page 138 for more information about pools properties.

`pset.min` Property and `pset.max` Property Constraints

These two properties place limits on the number of processors that can be allocated to a processor set, both minimum and maximum. See [Table 12-1](#) for more details about these properties.

Within these constraints, a resource partition’s resources are available to be allocated to other resource partitions in the same Solaris instance. Access to the resource is obtained by binding to a pool that is associated with the resource set. Binding is performed at login or manually by an administrator who has the `PRIV_SYS_RES_CONFIG` privilege.

`cpu.pinned` Property Constraint

The `cpu-pinned` property indicates that a particular CPU should not be moved by DRP from the processor set in which it is located. You can set this `libpool` property to maximize cache utilization for a particular application that is executing within a processor set.

See [Table 12-1](#) for more details about this property.

Configuration Objectives

Objectives are specified similarly to constraints. The full set of objectives is documented in [Table 12-1](#).

There are two categories of objectives.

- Workload dependent A workload-dependent objective is an objective that will vary according to the nature of the workload running on the system. An example is the `utilization` objective. The utilization figure for a resource set will vary according to the nature of the workload that is active in the set.
- Workload independent A workload-independent objective is an objective that does not vary according to the nature of the workload running on the system. An example is the `CPU locality` objective. The evaluated measure of locality for a resource set does not vary with the nature of the workload that is active in the set.

You can define three types of objectives.

Name	Valid Elements	Operators	Values
<code>wt-load</code>	<code>system</code>	N/A	N/A
<code>locality</code>	<code>pset</code>	N/A	<code>loose tight none</code>
<code>utilization</code>	<code>pset</code>	<code>< > ~</code>	0-100%

Objectives are stored in property strings in the `libpool` configuration. The property names are as follows:

- `system.poold.objectives`
- `pset.poold.objectives`

Objectives have the following syntax:

- `objectives = objective [; objective] *`
- `objective = [n:] keyword [op] [value]`

All objectives take an optional importance prefix. The importance acts as a multiplier for the objective and thus increases the significance of its contribution to the objective function evaluation. The range is from 0 to `INT64_MAX` (9223372036854775807). If not specified, the default importance value is 1.

Some element types support more than one type of objective. An example is `pset`. You can specify multiple objective types for these elements. You can also specify multiple utilization objectives on a single `pset` element.

See [“How to Define Configuration Objectives”](#) on page 165 for usage examples.

wt-load Objective

The `wt-load` objective favors configurations that match resource allocations to resource utilizations. A resource set that uses more resources will be given more resources when this objective is active.

Use this objective when you are satisfied with the constraints you have established using the minimum and maximum properties, and you would like the daemon to manipulate resources freely within those constraints.

The `locality` Objective

The `locality` objective influences the impact that locality, as measured by locality group (`lgroup`) data, has upon the selected configuration. An alternate definition for locality is latency. An `lgroup` describes CPU and memory resources. The `lgroup` is used by the Solaris system to determine the distance between resources, using time as the measurement. For more information on the locality group abstraction, see “Locality Groups Overview” in *Programming Interfaces Guide*.

This objective can take one of the following three values:

- `tight` If set, configurations that maximize resource locality are favored.
- `loose` If set, configurations that minimize resource locality are favored.
- `none` If set, the favorability of a configuration is not influenced by resource locality. This is the default value for the `locality` objective.

In general, the `locality` objective should be set to `tight`. However, to maximize memory bandwidth or to minimize the impact of DR operations on a resource set, you could set this objective to `loose` or keep it at the default setting of `none`.

utilization Objective

The `utilization` objective favors configurations that allocate resources to partitions that are not meeting the specified utilization objective.

This objective is specified by using operators and values. The operators are as follows:

- < The “less than” operator indicates that the specified value represents a maximum target value.
- > The “greater than” operator indicates that the specified value represents a minimum target value.
- ~ The “about” operator indicates that the specified value is a target value about which some fluctuation is acceptable.

A pset can only have one utilization objective set for each type of operator.

- If the ~ operator is set, then the < and > operators cannot be set.
- If the < and > operators are set, then the ~ operator cannot be set. Note that the settings of the < operator and the > operator cannot contradict each other.

You can set both a < and a > operator together to create a range. The values will be validated to make sure that they do not overlap.

Configuration Objectives Example

In the following example, `poold` is to assess these objectives for the pset:

- The `utilization` should be kept between 30 percent and 80 percent.
- The `locality` should be maximized for the processor set.
- The objectives should take the default importance of 1.

EXAMPLE 12-1 `poold` Objectives Example

```
pset.poold.objectives "utilization > 30; utilization < 80;  
locality tight"
```

See [“How to Define Configuration Objectives”](#) on page 165 for additional usage examples.

`poold` Properties

There are four categories of properties:

- Configuration
- Constraint
- Objective
- Objective Parameter

TABLE 12-1 Defined Property Names

Property Name	Type	Category	Description
<code>system.poold.log-level</code>	string	Configuration	Logging level
<code>system.poold.log-location</code>	string	Configuration	Logging location
<code>system.poold.monitor-interval</code>	unsigned int	Configuration	Monitoring sample interval
<code>system.poold.history-file</code>	string	Configuration	Decision history location

TABLE 12-1 Defined Property Names (Continued)

Property Name	Type	Category	Description
<code>system.pool.d.pid</code>	int	Configuration	<code>pool.d</code> PID
<code>pset.max</code>	unsigned int	Constraint	Maximum number of CPUs for this processor set
<code>pset.min</code>	unsigned int	Constraint	Minimum number of CPUs for this processor set
<code>cpu.pinned</code>	boolean	Constraint	CPUs pinned to this processor set
<code>system.pool.d.objectives</code>	string	Objective	Formatted string following <code>pool.d</code> 's objective expression syntax
<code>pset.pool.d.objectives</code>	string	Objective	Formatted string following <code>pool.d</code> 's expression syntax
<code>pool.importance</code>	unsigned int	Objective parameter	User-assigned importance

`pool.d` Features That Can Be Configured

You can configure these aspects of the daemon's behavior.

- Monitoring interval
- Logging level
- Logging location

These options are specified in the pools configuration. You can also control the logging level from the command line by invoking `pool.d`.

`pool.d` Monitoring Interval

Use the property name `system.pool.d.monitor-interval` to specify a value in milliseconds.

poold Logging Information

Three categories of information are provided through logging. These categories are identified in the logs:

- Configuration
- Monitoring
- Optimization

Use the property name `system.poold.log-level` to specify the logging parameter. If this property is not specified, the default logging level is `NOTICE`. The parameter levels are hierarchical. Setting a log level of `DEBUG` will cause `poold` to log all defined messages. The `INFO` level provides a useful balance of information for most administrators.

At the command line, you can use the `poold` command with the `-l` option and a parameter to specify the level of logging information generated.

The following parameters are available:

- ALERT
- CRIT
- ERR
- WARNING
- NOTICE
- INFO
- DEBUG

The parameter levels map directly onto their `syslog` equivalents. See [“Logging Location” on page 150](#) for more information about using `syslog`.

For more information about how to configure `poold` logging, see [“How to Set the poold Logging Level” on page 168](#).

Configuration Information Logging

The following types of messages can be generated:

ALERT	Problems accessing the <code>libpool</code> configuration, or some other fundamental, unanticipated failure of the <code>libpool</code> facility. Causes the daemon to exit and requires immediate administrative attention.
CRIT	Problems due to unanticipated failures. Causes the daemon to exit and requires immediate administrative attention.
ERR	Problems with the user-specified parameters that control operation, such as unresolvable, conflicting utilization objectives for a resource set. Requires administrative intervention to correct the objectives. <code>poold</code> attempts to take remedial action by ignoring conflicting objectives, but some errors will cause the daemon to exit.

- WARNING Warnings related to the setting of configuration parameters that, while technically correct, might not be suitable for the given execution environment. An example is marking all CPU resources as pinned, which means that `pool` cannot move CPU resources between processor sets.
- DEBUG Messages containing the detailed information that is needed when debugging configuration processing. This information is not generally used by administrators.

Monitoring Information Logging

The following types of messages can be generated:

- CRIT Problems due to unanticipated monitoring failures. Causes the daemon to exit and requires immediate administrative attention.
- ERR Problems due to unanticipated monitoring error. Could require administrative intervention to correct.
- NOTICE Messages about resource control region transitions.
- INFO Messages about resource utilization statistics.
- DEBUG Messages containing the detailed information that is needed when debugging monitoring processing. This information is not generally used by administrators.

Optimization Information Logging

The following types of messages can be generated:

- WARNING Messages could be displayed regarding problems making optimal decisions. Examples could include resource sets that are too narrowly constrained by their minimum and maximum values or by the number of pinned components.

Messages could be displayed about problems performing an optimal reallocation due to unforeseen limitations. Examples could include removing the last processor from a processor set which contains a bound resource consumer.
- NOTICE Messages about usable configurations or configurations that will not be implemented due to overriding decision histories could be displayed.
- INFO Messages about alternate configurations considered could be displayed.

DEBUG Messages containing the detailed information that is needed when debugging optimization processing. This information is not generally used by administrators.

Logging Location

The `system.pool.d.log-location` property is used to specify the location for `pool.d` logged output. You can specify a location of `SYSLOG` for `pool.d` output (see `syslog(3C)`).

If this property is not specified, the default location for `pool.d` logged output is `/var/log/pool/pool.d`.

When `pool.d` is invoked from the command line, this property is not used. Log entries are written to `stderr` on the invoking terminal.

Log Management With `logadm`

If `pool.d` is active, the `logadm.conf` file includes an entry to manage the default file `/var/log/pool/pool.d`. The entry is:

```
/var/log/pool/pool.d -N -s 512k
```

See the `logadm(1M)` and the `logadm.conf(4)` man pages.

How Dynamic Resource Allocation Works

This section explains the process and the factors that `pool.d` uses to dynamically allocate resources.

About Available Resources

Available resources are considered to be all of the resources that are available for use within the scope of the `pool.d` process. The scope of control is at most a single Solaris instance.

On a system that has zones enabled, the scope of an executing instance of `poold` is limited to the global zone.

Determining Available Resources

Resource pools encompass all of the system resources that are available for consumption by applications.

For a single executing Solaris instance, a resource of a single type, such as a CPU, must be allocated to a single partition. There can be one or more partitions for each type of resource. Each partition contains a unique set of resources.

For example, a machine with four CPUs and two processor sets can have the following setup:

```
pset 0: 0 1
```

```
pset 1: 2 3
```

where 0, 1, 2 and 3 after the colon represent CPU IDs. Note that the two processor sets account for all four CPUs.

The same machine cannot have the following setup:

```
pset 0: 0 1
```

```
pset 1: 1 2 3
```

It cannot have this setup because CPU 1 can appear in only one pset at a time.

Resources cannot be accessed from any partition other than the partition to which they belong.

To discover the available resources, `poold` interrogates the active pools configuration to find partitions. All resources within all partitions are summed to determine the total amount of available resources for each type of resource that is controlled.

This quantity of resources is the basic figure that `poold` uses in its operations. However, there are constraints upon this figure that limit the flexibility that `poold` has to make allocations. For information about available constraints, see [“Configuration Constraints” on page 143](#).

Identifying a Resource Shortage

The control scope for `pool` is defined as the set of available resources for which `pool` has primary responsibility for effective partitioning and management. However, other mechanisms that are allowed to manipulate resources within this control scope can still affect a configuration. If a partition should move out of control while `pool` is active, `pool` tries to restore control through the judicious manipulation of available resources. If `pool` cannot locate additional resources within its scope, then the daemon logs information about the resource shortage.

Determining Resource Utilization

`pool` typically spends the greatest amount of time observing the usage of the resources within its scope of control. This monitoring is performed to verify that workload-dependent objectives are being met.

For example, for processor sets, all measurements are made across all of the processors in a set. The resource utilization shows the proportion of time that the resource is in use over the sample interval. Resource utilization is displayed as a percentage from 0 to 100.

Identifying Control Violations

The directives described in [“Configuration Constraints and Objectives”](#) on page 142 are used to detect the approaching failure of a system to meet its objectives. These objectives are directly related to workload.

A partition that is not meeting user-configured objectives is a control violation. The two types of control violations are synchronous and asynchronous.

- A synchronous violation of an objective is detected by the daemon in the course of its workload monitoring.
- An asynchronous violation of an objective occurs independently of monitoring action by the daemon.

The following events cause asynchronous objective violations:

- Resources are added to or removed from a control scope.
- The control scope is reconfigured.
- The `pool` resource controller is restarted.

The contributions of objectives that are not related to workload are assumed to remain constant between evaluations of the objective function. Objectives that are not related to workload are only reassessed when a reevaluation is triggered through one of the asynchronous violations.

Determining Appropriate Remedial Action

When the resource controller determines that a resource consumer is short of resources, the initial response is that increasing the resources will improve performance.

Alternative configurations that meet the objectives specified in the configuration for the scope of control are examined and evaluated.

This process is refined over time as the results of shifting resources are monitored and each resource partition is evaluated for responsiveness. The decision history is consulted to eliminate reconfigurations that did not show improvements in attaining the objective function in the past. Other information, such as process names and quantities, are used to further evaluate the relevance of the historical data.

If the daemon cannot take corrective action, the condition is logged. For more information, see [“poolcd Logging Information” on page 148](#).

Using `poolstat` to Monitor the Pools Facility and Resource Utilization

The `poolstat` utility is used to monitor resource utilization when pools are enabled on your system. This utility iteratively examines all of the active pools on a system and reports statistics based on the selected output mode. The `poolstat` statistics enable you to determine which resource partitions are heavily utilized. You can analyze these statistics to make decisions about resource reallocation when the system is under pressure for resources.

The `poolstat` utility includes options that can be used to examine specific pools and report resource set-specific statistics.

If zones are implemented on your system and you use `poolstat` in a non-global zone, information about the resources associated with the zone’s pool is displayed.

For more information about the `poolstat` utility, see the `poolstat(1M)` man page. For `poolstat` task and usage information, see [“Using `poolstat` to Report Statistics for Pool-Related Resources” on page 173](#).

poolstat Output

In default output format, `poolstat` outputs a heading line and then displays a line for each pool. A pool line begins with the pool ID and the name of the pool, followed by a column of statistical data for the processor set attached to the pool. Resource sets attached to more than one pool are listed multiple times, once for each pool.

The column headings are as follows:

<code>id</code>	Pool ID.
<code>pool</code>	Pool name.
<code>rid</code>	Resource set ID.
<code>rset</code>	Resource set name.
<code>type</code>	Resource set type.
<code>min</code>	Minimum resource set size.
<code>max</code>	Maximum resource set size.
<code>size</code>	Current resource set size.
<code>used</code>	Measure of how much of the resource set is currently used.

This usage is calculated as the percentage of utilization of the resource set multiplied by the size of the resource set. If a resource set has been reconfigured during the last sampling interval, this value might be not reported. An unreported value appears as a hyphen (-).

<code>load</code>	Absolute representation of the load that is put on the resource set.
-------------------	--

For more information about this property, see the `libpool(3LIB)` man page.

You can specify the following in `poolstat` output:

- The order of the columns
- The headings that appear

Tuning poolstat Operation Intervals

You can customize the operations performed by `poolstat`. You can set the sampling interval for the report and specify the number of times that statistics are repeated:

<i>interval</i>	Tune the intervals for the periodic operations performed by <code>poolstat</code> . All intervals are specified in seconds.
<i>count</i>	Specify the number of times that the statistics are repeated. By default, <code>poolstat</code> reports statistics only once.

If *interval* and *count* are not specified, statistics are reported once. If *interval* is specified and *count* is not specified, then statistics are reported indefinitely.

Commands Used With the Resource Pools Facility

The commands described in the following table provide the primary administrative interface to the pools facility. For information on using these commands on a system that has zones enabled, see [“Resource Pools Used in Zones”](#) on page 135.

Man Page Reference	Description
<code>pooladm(1M)</code>	Enables or disables the pools facility on your system. Activates a particular configuration or removes the current configuration and returns associated resources to their default status. If run without options, <code>pooladm</code> prints out the current dynamic pools configuration.
<code>poolbind(1M)</code>	Enables the manual binding of projects, tasks, and processes to a resource pool.
<code>poolcfg(1M)</code>	Provides configuration operations on pools and sets. Configurations created using this tool are instantiated on a target host by using <code>pooladm</code> . If run with the <code>info</code> subcommand argument to the <code>-c</code> option, <code>poolcfg</code> displays information about the static configuration at <code>/etc/pooladm.conf</code> . If a file name argument is added, this command displays information about the static configuration held in the named file. For example, <code>poolcfg -c info /tmp/newconfig</code> displays information about the static configuration contained in the file <code>/tmp/newconfig</code> .
<code>poold(1M)</code>	The pools system daemon. The daemon uses system targets and observable statistics to preserve the system performance goals specified by the administrator. If unable to take corrective action when goals are not being met, <code>poold</code> logs the condition.
<code>poolstat(1M)</code>	Displays statistics for pool-related resources. Simplifies performance analysis and provides information that supports system administrators in resource partitioning and repartitioning tasks. Options are provided for examining specified pools and reporting resource set-specific statistics.

A library API is provided by `libpool` (see the `libpool(3LIB)` man page). The library can be used by programs to manipulate pool configurations.

Administering Dynamic Resource Pools (Tasks)

This chapter describes how to set up and administer resource pools on your system.

For background information about resource pools, see [Chapter 12](#).

Administering Dynamic Resource Pools (Task Map)

Task	Description	For Instructions
Enable the resource pools facility.	Activate the resource pools framework on your system.	“How to Enable Pools” on page 159
Disable the resource pools facility.	Disable the resource pools framework on your system.	“How to Disable Pools” on page 159
Create a static resource pools configuration.	Create a static configuration file that matches the current dynamic configuration. For more information, see “Resource Pools Framework” on page 137 .	“How to Create a Static Configuration” on page 160
Modify a resource pools configuration.	Revise a pools configuration on your system, for example, by creating additional pools.	“How to Modify a Configuration” on page 161

Task	Description	For Instructions
Associate a resource pool with a scheduling class.	Associate a pool with a scheduling class so that all processes bound to the pool use the specified scheduler.	“How to Associate a Pool With a Scheduling Class” on page 163
Define configuration objectives.	Specify objectives for <code>poold</code> to consider when taking corrective action. For more information on configuration objectives, see “poold Overview” on page 141.	“How to Define Configuration Objectives” on page 165
Set the logging level.	Specify the level of logging information that <code>poold</code> generates.	“How to Set the poold Logging Level” on page 168
Use a text file with the <code>poolcfg</code> command.	The <code>poolcfg</code> command can take input from a text file.	“How to Use Command Files With poolcfg” on page 168
Transfer resources in the kernel.	Transfer resources in the kernel. For example, transfer resources with specific IDs to a target set.	“Transferring Resources” on page 169
Activate a pools configuration.	Activate the configuration in the default configuration file.	“How to Activate a Pools Configuration” on page 170
Validate a pools configuration before you commit the configuration.	Validate a pools configuration to test what will happen when the validation occurs.	“How to Validate a Configuration Before Committing the Configuration” on page 170
Remove a pools configuration from your system.	All associated resources, such as processor sets, are returned to their default status.	“How to Remove a Pools Configuration” on page 170
Bind processes to a pool.	Manually associate a running process on your system with a resource pool.	“How to Bind Processes to a Pool” on page 171
Bind tasks or projects to a pool.	Associate tasks or projects with a resource pool.	“How to Bind Tasks or Projects to a Pool” on page 172
Bind new processes to a resource pool.	To automatically bind new processes in a project to a given pool, add an attribute to each entry in the <code>project</code> database.	“How to Use project Attributes to Bind New Processes to a Pool” on page 172

Task	Description	For Instructions
Use <code>project</code> attributes to bind a process to a different pool.	Modify the pool binding for new processes that are started.	“How to Use <code>project</code> Attributes to Bind a Process to a Different Pool” on page 173
Use the <code>poolstat</code> utility to produce reports.	Produce multiple reports at specified intervals.	“Producing Multiple Reports at Specific Intervals” on page 173
Report resource set statistics.	Use the <code>poolstat</code> utility to report statistics for a <code>pset</code> resource set.	“Reporting Resource Set Statistics” on page 174

Enabling and Disabling the Pools Facility

You can use the `pooladm` command described in the `pooladm(1M)` man page to perform the following tasks:

- Enable the pools facility so that pools can be manipulated
- Disable the pools facility so that pools cannot be manipulated

▼ How to Enable Pools

1. **Become superuser, or assume a role that includes the Process Management profile.**

The System Administrator role includes the Process Management profile. For information on how to create the role and assign the role to a user, see *“Managing RBAC (Task Map)”* in *System Administration Guide: Security Services*.

2. **Enable the pools facility.**

```
# pooladm -e
```

▼ How to Disable Pools

1. **Become superuser, or assume a role that includes the Process Management profile.**

The System Administrator role includes the Process Management profile. For information on how to create the role and assign the role to a user, see *“Managing RBAC (Task Map)”* in *System Administration Guide: Basic Administration*.

2. Disable the pools facility.

```
# pooladm -d
```

Configuring Pools

▼ How to Create a Static Configuration

Enable pools on your system, then use the `-s` option to `/usr/sbin/pooladm` to create a static configuration file that matches the current dynamic configuration. Unless a different file name is specified, the default location `/etc/pooladm.conf` is used.

Commit your configuration using the `pooladm` command with the `-c` option. Then, use the `pooladm` command with the `-s` option to update the static configuration to match the state of the dynamic configuration.

1. Become superuser, or assume a role that includes the Process Management profile.

The System Administrator role includes the Process Management profile. For information on how to create the role and assign the role to a user, see “Managing RBAC (Task Map)” in *System Administration Guide: Security Services*.

2. Update the static configuration file to match the current dynamic configuration.

```
# pooladm -s
```

3. View the contents of the configuration file in readable form.

Note that the configuration contains default elements created by the system.

```
# poolcfg -c info
system tester
  string system.comment
  int    system.version 1
  boolean system.bind-default true
  int    system.poold.pid 177916

  pool pool_default
    int    pool.sys_id 0
    boolean pool.active true
    boolean pool.default true
    int    pool.importance 1
    string pool.comment
    pset   pset_default
```



```

pset pset_default
    int      pset.sys_id -1
    boolean  pset.default true
    uint     pset.min 1
    uint     pset.max 65536
    string   pset.units population
    uint     pset.load 10
    uint     pset.size 4
    string   pset.comment
    boolean  testnullchanged true

    cpu

        int      cpu.sys_id 3
        string   cpu.comment
        string   cpu.status on-line

    cpu

        int      cpu.sys_id 2
        string   cpu.comment
        string   cpu.status on-line

    cpu

        int      cpu.sys_id 1
        string   cpu.comment
        string   cpu.status on-line

    cpu

        int      cpu.sys_id 0
        string   cpu.comment
        string   cpu.status on-line

```

4. Commit the configuration at `/etc/pooladm.conf`.

```
# pooladm -c
```

5. (Optional) To copy the dynamic configuration to a static configuration file called `/tmp/backup`, type the following:

```
# pooladm -s /tmp/backup
```

Note – The new functionality `pooladm -s` is preferred over the previous functionality `poolcfg -c discover` for creating a new configuration that matches the dynamic configuration.

▼ How to Modify a Configuration

To enhance your configuration, create a processor set named `pset_batch` and a pool named `pool_batch`. Then join the pool and the processor set with an association.

Note that you must quote subcommand arguments that contain white space.

1. Become superuser, or assume a role that includes the Process Management profile.

The System Administrator role includes the Process Management profile. For information on how to create the role and assign the role to a user, see “Managing RBAC” in *System Administration Guide: Security Services*.

2. Create processor set `pset_batch`.

```
# poolcfg -c 'create pset pset_batch (uint pset.min = 2; uint pset.max = 10)'
```

3. Create pool `pool_batch`.

```
# poolcfg -c 'create pool pool_batch'
```

4. Join the pool and the processor set with an association.

```
# poolcfg -c 'associate pool pool_batch (pset pset_batch)'
```

5. Display the edited configuration.

```
# poolcfg -c info
system tester
  string system.comment kernel state
  int    system.version 1
  boolean system.bind-default true
  int    system.poold.pid 177916

  pool pool_default
    int    pool.sys_id 0
    boolean pool.active true
    boolean pool.default true
    int    pool.importance 1
    string pool.comment
    pset   pset_default

  pset pset_default
    int    pset.sys_id -1
    boolean pset.default true
    uint   pset.min 1
    uint   pset.max 65536
    string pset.units population
    uint   pset.load 10
    uint   pset.size 4
    string pset.comment
    boolean testnullchanged true

  cpu
    int    cpu.sys_id 3
    string cpu.comment
    string cpu.status on-line

  cpu
    int    cpu.sys_id 2
    string cpu.comment
    string cpu.status on-line
```

```

cpu
    int    cpu.sys_id 1
    string cpu.comment
    string cpu.status on-line

cpu
    int    cpu.sys_id 0
    string cpu.comment
    string cpu.status on-line

pool pool_batch
    boolean pool.default false
    boolean pool.active true
    int pool.importance 1
    string pool.comment
    pset pset_batch

pset pset_batch
    int pset.sys_id -2
    string pset.units population
    boolean pset.default true
    uint pset.max 10
    uint pset.min 2
    string pset.comment
    boolean pset.escapable false
    uint pset.load 0
    uint pset.size 0

cpu
    int    cpu.sys_id 5
    string cpu.comment
    string cpu.status on-line

cpu
    int    cpu.sys_id 4
    string cpu.comment
    string cpu.status on-line

```

6. Commit the configuration at `/etc/pooladm.conf`.

```
# pooladm -c
```

7. (Optional) To copy the dynamic configuration to a static configuration file named `/tmp/backup`, type the following:

```
# pooladm -s /tmp/backup
```

▼ How to Associate a Pool With a Scheduling Class

You can associate a pool with a scheduling class so that all processes bound to the pool use this scheduler. To do this, set the `pool.scheduler` property to the name of the scheduler. This example associates the pool `pool_batch` with the fair share scheduler (FSS).

1. Become superuser, or assume a role that includes the Process Management profile.

The System Administrator role includes the Process Management profile. For information on how to create the role and assign the role to a user, see “Managing RBAC (Task Map)” in *System Administration Guide: Security Services*.

2. Modify pool pool_batch to be associated with the FSS.

```
# poolcfg -c 'modify pool pool_batch (string pool.scheduler="FSS")'
```

3. Display the edited configuration.

```
# poolcfg -c info
system tester
  string system.comment
  int    system.version 1
  boolean system.bind-default true
  int    system.poold.pid 177916

  pool pool_default
    int    pool.sys_id 0
    boolean pool.active true
    boolean pool.default true
    int    pool.importance 1
    string pool.comment
    pset   pset_default

  pset pset_default
    int    pset.sys_id -1
    boolean pset.default true
    uint   pset.min 1
    uint   pset.max 65536
    string pset.units population
    uint   pset.load 10
    uint   pset.size 4
    string pset.comment
    boolean testnullchanged true

  cpu
    int    cpu.sys_id 3
    string cpu.comment
    string cpu.status on-line

  cpu
    int    cpu.sys_id 2
    string cpu.comment
    string cpu.status on-line

  cpu
    int    cpu.sys_id 1
    string cpu.comment
    string cpu.status on-line

  cpu
```

```

        int    cpu.sys_id 0
        string cpu.comment
        string cpu.status on-line

pool pool_batch
    boolean pool.default false
    boolean pool.active true
    int pool.importance 1
    string pool.comment
    string pool.scheduler FSS
    pset batch

pset pset_batch
    int pset.sys_id -2
    string pset.units population
    boolean pset.default true
    uint pset.max 10
    uint pset.min 2
    string pset.comment
    boolean pset.escapable false
    uint pset.load 0
    uint pset.size 0

cpu
    int    cpu.sys_id 5
    string cpu.comment
    string cpu.status on-line

cpu
    int    cpu.sys_id 4
    string cpu.comment
    string cpu.status on-line

```

4. Commit the configuration at `/etc/pooladm.conf`:

```
# pooladm -c
```

5. (Optional) To copy the dynamic configuration to a static configuration file called `/tmp/backup`, type the following:

```
# pooladm -s /tmp/backup
```

▼ How to Define Configuration Objectives

You can specify objectives for `pool` to consider when taking corrective action.

In the following procedure, the `wt-load` objective is being set so that `pool` tries to match resource allocation to resource utilization. The `locality` objective is disabled to assist in achieving this configuration goal.

- 1. Become superuser, or assume a role that includes the Process Management profile.**

The System Administrator role includes the Process Management profile. For information on how to create the role and assign the role to a user, see “Managing RBAC” in *System Administration Guide: Security Services*.

2. Modify system tester to favor the wt-load objective.

```
# poolcfg -c 'modify system tester (string system.poold.objectives="wt-load")'
```

3. Disable the locality objective for the default processor set.

```
# poolcfg -c 'modify pset pset_default (string pset.poold.objectives="locality none")'
```

4. Disable the locality objective for the pset_batch processor set.

```
# poolcfg -c 'modify pset pset_batch (string pset.poold.objectives="locality none")'
```

5. Display the edited configuration.

```
# poolcfg -c info
system tester
  string system.comment
  int system.version 1
  boolean system.bind-default true
  int system.poold.pid 177916
  string system.poold.objectives wt-load

pool pool_default
  int pool.sys_id 0
  boolean pool.active true
  boolean pool.default true
  int pool.importance 1
  string pool.comment
  pset pset_default

pset pset_default
  int pset.sys_id -1
  boolean pset.default true
  uint pset.min 1
  uint pset.max 65536
  string pset.units population
  uint pset.load 10
  uint pset.size 4
  string pset.comment
  boolean testnullchanged true
  string pset.poold.objectives locality none

cpu
  int cpu.sys_id 3
  string cpu.comment
  string cpu.status on-line

cpu
  int cpu.sys_id 2
  string cpu.comment
  string cpu.status on-line
```

```

cpu
    int    cpu.sys_id 1
    string cpu.comment
    string cpu.status on-line

cpu
    int    cpu.sys_id 0
    string cpu.comment
    string cpu.status on-line

pool pool_batch
    boolean pool.default false
    boolean pool.active true
    int pool.importance 1
    string pool.comment
    string pool.scheduler FSS
    pset batch

pset pset_batch
    int pset.sys_id -2
    string pset.units population
    boolean pset.default true
    uint pset.max 10
    uint pset.min 2
    string pset.comment
    boolean pset.escapable false
    uint pset.load 0
    uint pset.size 0
    string pset.poolid.objectives locality none

cpu
    int    cpu.sys_id 5
    string cpu.comment
    string cpu.status on-line

cpu
    int    cpu.sys_id 4
    string cpu.comment
    string cpu.status on-line

```

6. Commit the configuration at `/etc/pooladm.conf`.

```
# pooladm -c
```

7. (Optional) To copy the dynamic configuration to a static configuration file called `/tmp/backup`, type the following:

```
# pooladm -s /tmp/backup
```

▼ How to Set the `pool`d Logging Level

To specify the level of logging information that `pool`d generates, set the `system.pool.d.log-level` property in the `pool`d configuration. The `pool`d configuration is held in the `libpool` configuration. For information, see “[poold Logging Information” on page 148 and the `poolcfg\(1M\)` and `libpool\(3LIB\)` man pages.](#)

You can also use the `pool`d command at the command line to specify the level of logging information that `pool`d generates.

1. **Become superuser, or assume a role that includes the Process Management profile.**

The System Administrator role includes the Process Management profile. For information on how to create the role and assign the role to a user, see “Managing RBAC (Task Map)” in *System Administration Guide: Security Services*.

2. **Set the logging level by using the `pool`d command with the `-l` option and a parameter, for example, `INFO`.**

```
# /usr/lib/pool/poold -l INFO
```

For information about available parameters, see “[poold Logging Information” on page 148. The default logging level is `NOTICE`.](#)

▼ How to Use Command Files With `pool`cfg

The `pool`cfg command with the `-f` option can take input from a text file that contains `pool`cfg subcommand arguments to the `-c` option. This method is appropriate when you want a set of operations to be performed. When processing multiple commands, the configuration is only updated if all of the commands succeed. For large or complex configurations, this technique can be more useful than per-subcommand invocations.

Note that in command files, the `#` character acts as a comment mark for the rest of the line.

1. **Create the input file `poolcmds.txt`.**

```
$ cat > poolcmds.txt
create system tester
create pset pset_batch (uint pset.min = 2; uint pset.max = 10)
create pool pool_batch
associate pool pool_batch (pset pset_batch)
```

2. **Become superuser, or assume a role that includes the Process Management profile.**

The System Administrator role includes the Process Management profile. For information on how to create the role and assign the role to a user, see “Managing

RBAC” in *System Administration Guide: Security Services*.

3. Execute the command:

```
# /usr/sbin/poolcfg -f poolcmds.txt
```

Transferring Resources

Use the `transfer` subcommand argument to the `-c` option of `poolcfg` with the `-d` option to transfer resources in the kernel. The `-d` option specifies that the command operate directly on the kernel and not take input from a file. The following procedure moves two CPUs from processor set `pset1` to processor set `pset2` in the kernel.

1. Become superuser, or assume a role that includes the Process Management profile.

The System Administrator role includes the Process Management profile. For information on how to create the role and assign the role to a user, see “Managing RBAC” in *System Administration Guide: Security Services*.

2. Move two CPUs from `pset1` to `pset2`.

```
# poolcfg -dc 'transfer 2 from pset pset1 to pset2'
```

The `from` and `to` subclauses can be used in any order. Only one `to` and `from` subclause is supported per command.

If specific known IDs of a resource type are to be transferred, an alternative syntax is provided. For example, the following command assigns two CPUs with IDs 0 and 2 to the `pset_large` processor set:

```
# poolcfg -dc "transfer to pset pset_large (cpu 0; cpu 2)"
```

If a transfer fails because there are not enough resources to match the request or because the specified IDs cannot be located, the system displays an error message.

Activating and Removing Pool Configurations

Use the `pooladm` command to make a particular pool configuration active or to remove the currently active pool configuration. See the `pooladm(1M)` man page for more information about this command.

▼ How to Activate a Pools Configuration

To activate the configuration in the default configuration file, `/etc/pooladm.conf`, invoke `pooladm` with the `-c` option, “commit configuration.”

1. **Become superuser, or assume a role that includes the Process Management profile.**

The System Administrator role includes the Process Management profile. For information on how to create the role and assign the role to a user, see “Managing RBAC” in *System Administration Guide: Security Services*.

2. **Commit the configuration at `/etc/pooladm.conf`.**

```
# pooladm -c
```

3. **(Optional) Copy the dynamic configuration to a static configuration file, for example, `/tmp/backup`.**

```
# pooladm -s /tmp/backup
```

▼ How to Validate a Configuration Before Committing the Configuration

You can use the `-n` option with the `-c` option to test what will happen when the validation occurs. The configuration will not actually be committed.

The following command attempts to validate the configuration contained at `/home/admin/newconfig`. Any error conditions encountered are displayed, but the configuration itself is not modified.

1. **Become superuser, or assume a role that includes the Process Management profile.**

The System Administrator role includes the Process Management profile. For information on how to create the role and assign the role to a user, see “Managing RBAC (Task Map)” in *System Administration Guide: Security Services*.

2. **Test the validity of the configuration before committing it.**

```
# pooladm -n -c /home/admin/newconfig
```

▼ How to Remove a Pools Configuration

To remove the current active configuration and return all associated resources, such as processor sets, to their default status, use the `-x` option for “remove configuration.”

1. **Become superuser, or assume a role that includes the Process Management profile.**

The System Administrator role includes the Process Management profile. For information on how to create the role and assign the role to a user, see “Managing RBAC (Task Map)” in *System Administration Guide: Security Services*.

2. Remove the current active configuration.

```
# pooladm -x
```

The `-x` option to `pooladm` removes all user-defined elements from the dynamic configuration. All resources revert to their default states, and all pool bindings are replaced with a binding to the default pool.

Note – Mixing scheduling classes within one processor set can lead to unpredictable results. If the use of `pooladm -x` results in mixed scheduling classes within one processor set, use the `priocntl` command to move running processes into a different scheduling class. An example is provided in “How to Manually Move Processes From the TS Class Into the FSS Class” on page 109. You can also see the `priocntl(1)` man page.

Binding to a Pool

You can bind a running process to a pool in two ways:

- You can use the `poolbind` command described in `poolbind(1M)` command to bind a specific process to a named resource pool.
- You can use the `project.pool` attribute in the `project` database to identify the pool binding for a new login session or a task that is launched through the `newtask` command. See the `newtask(1)` and `project(4)` man pages.

▼ How to Bind Processes to a Pool

The following procedure uses `poolbind` with the `-p` option to manually bind a process (in this case, the current shell) to a pool named `ohare`.

1. Become superuser, or assume a role that includes the Process Management profile.

The System Administrator role includes the Process Management profile. For information on how to create the role and assign the role to a user, see “Managing RBAC (Task Map)” in *System Administration Guide: Security Services*.

2. Manually bind a process to a pool:

```
# poolbind -p ohare $$
```

3. Verify the pool binding for the process by using `poolbind` with the `-q` option.

```
$ poolbind -q $$  
155509 ohare
```

The system displays the process ID and the pool binding.

▼ How to Bind Tasks or Projects to a Pool

To bind tasks or projects to a pool, use the `poolbind` command with the `-i` option. The following example binds all processes in the `airmiles` project to the `laguardia` pool.

1. Become superuser, or assume a role that includes the Process Management profile.

The System Administrator role includes the Process Management profile. For information on how to create the role and assign the role to a user, see “Managing RBAC (Task Map)” in *System Administration Guide: Security Services*.

2. Bind all processes in the `airmiles` project to the `laguardia` pool.

```
# poolbind -i project -p laguardia airmiles
```

▼ How to Use project Attributes to Bind New Processes to a Pool

1. Become superuser, or assume a role that includes the Process Management profile.

The System Administrator role includes the Process Management profile. For information on how to create the role and assign the role to a user, see “Managing RBAC (Task Map)” in *System Administration Guide: Security Services*.

2. To automatically bind new processes in a project to a given pool, add the `project.pool` attribute to each entry in the `project` database.

For example, assume you have a configuration with two pools that are named `studio` and `backstage`. The `/etc/project` file would have the following contents:

```
user.paul:1024:::project.pool=studio  
user.george:1024:::project.pool=studio  
user.ringo:1024:::project.pool=backstage  
passes:1027::paul::project.pool=backstage
```

With this configuration, processes that are started by user `paul` are bound by default to the `studio` pool.

▼ How to Use `project` Attributes to Bind a Process to a Different Pool

Using the previous configuration, user `paul` can modify the pool binding for processes he starts. `paul` can use `newtask` to bind work to the `backstage` pool as well, by launching in the `passes` project.

1. Launch a process in the `passes` project.

```
$ newtask -l -p passes
```

2. Use the `poolbind` command with the `-q` option to verify the pool binding for the process.

```
$ poolbind -q $$  
6384 pool backstage
```

The system displays the process ID and the pool binding.

Using `poolstat` to Report Statistics for Pool-Related Resources

The `poolstat` command is used to display statistics for pool-related resources. See [“Using `poolstat` to Monitor the Pools Facility and Resource Utilization”](#) on page 153 and the `poolstat(1M)` man page for more information.

The following subsections use examples to illustrate how to produce reports for specific purposes.

Displaying Default `poolstat` Output

Typing `poolstat` without arguments outputs a header line and a line of information for each pool. The information line shows the pool ID, the name of the pool, and resource statistics for the processor set attached to the pool.

```
machine% poolstat  
  
           pset  
id pool      size used load  
0 pool_default 4  3.6  6.2  
1 pool_sales   4  3.3  8.4
```

Producing Multiple Reports at Specific Intervals

The following command produces three reports at 5-second sampling intervals.

```

machine% poolstat 5 3

      pset
id pool      size used load
46 pool_sales 2  1.2  8.3
 0 pool_default 2  0.4  5.2

      pset
id pool      size used load
46 pool_sales 2  1.4  8.4
 0 pool_default 2  1.9  2.0

      pset
id pool      size used load
46 pool_sales 2  1.1  8.0
 0 pool_default 2  0.3  5.0

```

Reporting Resource Set Statistics

The following example uses the `poolstat` command with the `-r` option to report statistics for the processor set resource set. Note that the resource set `pset_default` is attached to more than one pool, so this processor set is listed once for each pool membership.

```

machine% poolstat -r pset
      id pool      type rid rset      min  max size used load
      0 pool_default pset -1 pset_default  1  65K  2  1.2  8.3
      6 pool_sales   pset  1 pset_sales   1  65K  2  1.2  8.3
      2 pool_other   pset -1 pset_default  1  10K  2  0.4  5.2

```

Resource Management Configuration Example

This chapter reviews the resource management framework and describes a hypothetical server consolidation project.

The following topics are covered in this chapter:

- “Configuration to Be Consolidated” on page 175
- “Consolidation Configuration” on page 176
- “Creating the Configuration” on page 176
- “Viewing the Configuration” on page 178

Configuration to Be Consolidated

In this example, five applications are being consolidated onto a single system. The target applications have resource requirements that vary, different user populations, and different architectures. Currently, each application exists on a dedicated server that is designed to meet the requirements of the application. The applications and their characteristics are identified in the following table.

Application Description	Characteristics
Application server	Exhibits negative scalability beyond 2 CPUs
Database instance for application server	Heavy transaction processing
Application server in test and development environment	GUI-based, with untested code execution
Transaction processing server	Primary concern is response time

Application Description	Characteristics
Standalone database instance	Processes a large number of transactions and serves multiple time zones

Consolidation Configuration

The following configuration is used to consolidate the applications onto a single system.

- The application server has a two-CPU processor set.
- The database instance for the application server and the standalone database instance are consolidated onto a single processor set of at least four CPUs. The standalone database instance is guaranteed 75 percent of that resource.
- The test and development application server requires the IA scheduling class to ensure UI responsiveness. Memory limitations are imposed to lessen the effects of bad code builds.
- The transaction processing server is assigned a dedicated processor set of at least two CPUs, to minimize response latency.

This configuration covers known applications that are executing and consuming processor cycles in each resource set. Thus, constraints can be established that allow the processor resource to be transferred to sets where the resource is required.

- The `wt-load` objective is set to allow resource sets that are highly utilized to receive greater resource allocations than sets that have low utilization.
- The `locality` objective is set to `tight`, which is used to maximize processor locality.

An additional constraint to prevent utilization from exceeding 80 percent of any resource set is also applied. This constraint ensures that applications get access to the resources they require. Moreover, for the transaction processor set, the objective of maintaining utilization below 80 percent is twice as important as any other objectives that are specified. This importance will be defined in the configuration.

Creating the Configuration

Edit the `/etc/project` database file. Add entries to implement the required resource controls and to map users to resource pools, then view the file.


```
# 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 development::staff:project.pool=dev_pool;
process.max-address-space=(privileged,536870912,deny)    keep with previous line
user.tp_engine:2004:Transaction Engine:::project.pool=tp_pool
user.geo_db:2005:EDI DB:::project.pool=db_pool;project.cpu-shares=(privileged,3,deny)
.
.
.
```

Note – The development team has to execute tasks in the development project because access for this project is based on a user’s group ID (GID).

Create an input file named `pool.host`, which will be used to configure the required resource pools. View the file.

```
# cat pool.host
create system host
create pset dev_pset (uint pset.min = 0; uint pset.max = 2)
create pset tp_pset (uint pset.min = 2; uint pset.max=8)
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 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)
modify system tester (string system.poold.objectives="wtload")
modify pset dev_pset (string pset.poold.objectives="locality tight; utilization < 80")
modify pset tp_pset (string pset.poold.objectives="locality tight; 2: utilization < 80")
modify pset db_pset (string pset.poold.objectives="locality tight;utilization < 80")
modify pset app_pset (string pset.poold.objectives="locality tight; utilization < 80")
```

Update the configuration using the `pool.host` input file.

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

Make the configuration active.

```
# pooladm -c
```

The framework is now functional on the system.

Viewing the Configuration

To view the framework configuration, which also contains default elements created by the system, type:

```
# pooladm
system host
  string system.comment
  int system.version 1
  boolean system.bind-default true
  int system.poolid.pid 177916
  string system.poolid.objectives wtlload

pool dev_pool
  int pool.sys_id 125
  boolean pool.default false
  boolean pool.active true
  int pool.importance 1
  string pool.comment
  string pool.scheduler IA
  pset dev_pset

pool appserver_pool
  int pool.sys_id 124
  boolean pool.default false
  boolean pool.active true
  int pool.importance 1
  string pool.comment
  string pool.scheduler TS
  pset app_pset

pool db_pool
  int pool.sys_id 123
  boolean pool.default false
  boolean pool.active true
  int pool.importance 1
  string pool.comment
  string pool.scheduler FSS
  pset db_pset

pool tp_pool
  int pool.sys_id 122
  boolean pool.default false
  boolean pool.active true
  int pool.importance 1
  string pool.comment
  string pool.scheduler TS
  pset tp_pset

pool pool_default
  int pool.sys_id 0
  boolean pool.default true
```

```

        boolean pool.active true
        int     pool.importance 1
        string  pool.comment
        string  pool.scheduler TS
        pset   pset_default

pset dev_pset
    int     pset.sys_id 4
    string  pset.units population
    boolean pset.default false
    uint    pset.min 0
    uint    pset.max 2
    string  pset.comment
    boolean pset.escapable false
    uint    pset.load 0
    uint    pset.size 0
    string  pset.poold.objectives locality tight; utilization < 80

pset tp_pset
    int     pset.sys_id 3
    string  pset.units population
    boolean pset.default false
    uint    pset.min 2
    uint    pset.max 8
    string  pset.comment
    boolean pset.escapable false
    uint    pset.load 0
    uint    pset.size 0
    string  pset.poold.objectives locality tight; 2: utilization < 80

    cpu
        int     cpu.sys_id 1
        string  cpu.comment
        string  cpu.status on-line

    cpu
        int     cpu.sys_id 2
        string  cpu.comment
        string  cpu.status on-line

pset db_pset
    int     pset.sys_id 2
    string  pset.units population
    boolean pset.default false
    uint    pset.min 4
    uint    pset.max 6
    string  pset.comment
    boolean pset.escapable false
    uint    pset.load 0
    uint    pset.size 0
    string  pset.poold.objectives locality tight; utilization < 80

    cpu
        int     cpu.sys_id 3
        string  cpu.comment

```

```

        string  cpu.status on-line

cpu
    int      cpu.sys_id 4
    string   cpu.comment
    string   cpu.status on-line

cpu
    int      cpu.sys_id 5
    string   cpu.comment
    string   cpu.status on-line

cpu
    int      cpu.sys_id 6
    string   cpu.comment
    string   cpu.status on-line

pset app_pset
    int      pset.sys_id 1
    string   pset.units population
    boolean  pset.default false
    uint     pset.min 1
    uint     pset.max 2
    string   pset.comment
    boolean  pset.escapable false
    uint     pset.load 0
    uint     pset.size 0
    string   pset.pool.default.objectives locality tight; utilization < 80
    cpu
        int      cpu.sys_id 7
        string   cpu.comment
        string   cpu.status on-line

pset pset_default
    int      pset.sys_id -1
    string   pset.units population
    boolean  pset.default true
    uint     pset.min 1
    uint     pset.max 4294967295
    string   pset.comment
    boolean  pset.escapable false
    uint     pset.load 0
    uint     pset.size 0

cpu
    int      cpu.sys_id 0
    string   cpu.comment
    string   cpu.status on-line

```

A graphic representation of the framework follows.

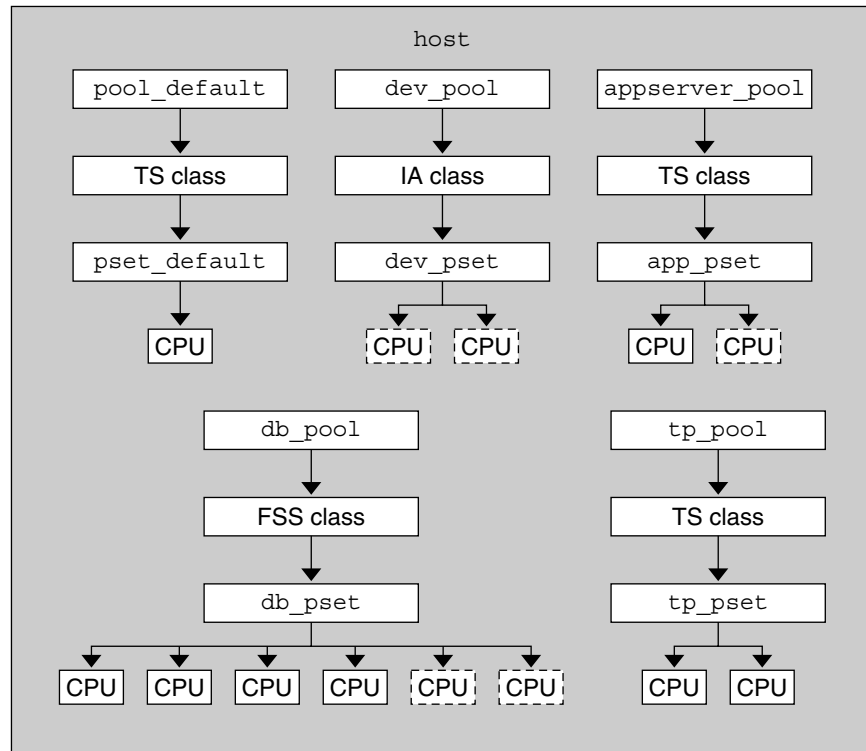


FIGURE 14-1 Server Consolidation Configuration

Note – In the pool `db_pool`, the standalone database instance is guaranteed 75 percent of the CPU resource.

Resource Control Functionality in the Solaris Management Console

This chapter describes the resource control and performance monitoring features in the Solaris Management Console. Only a subset of the resource management features can be controlled using the console.

You can use the console to monitor system performance and to enter the resource control values shown in [Table 15–1](#) for projects, tasks, and processes. The console provides a convenient, secure alternative to the command-line interface (CLI) for managing hundreds of configuration parameters that are spread across many systems. Each system is managed individually. The console’s graphical interface supports all experience levels.

The following topics are covered.

- “Using the Console (Task Map)” on page 184
- “Console Overview” on page 184
- “Management Scope” on page 184
- “Performance Tool” on page 185
- “Resource Controls Tab” on page 188
- “Console References” on page 191

Using the Console (Task Map)

Task	Description	For Instructions
Use the console	Start the Solaris Management Console in a local environment or in a name service or directory service environment. Note that the performance tool is not available in a name service environment.	“Starting the Solaris Management Console” in <i>System Administration Guide: Basic Administration</i> and “Using the Solaris Management Tools in a Name Service Environment (Task Map)” in <i>System Administration Guide: Basic Administration</i>
Monitor system performance	Access the Performance tool under System Status.	“How to Access the Performance Tool” on page 185
Add resource controls to projects	Access the Resource Controls tab under System Configuration.	“How to Access the Resource Controls Tab” on page 189

Console Overview

Resource management functionality is a component of the Solaris Management Console. The console is a container for GUI-based administrative tools that are stored in collections called toolboxes. For information on the console and how to use it, see [“Working With the Management Console \(Tasks\)”](#) in *System Administration Guide: Basic Administration*.

When you use the console and its tools, the main source of documentation is the online help system in the console itself. For a description of the documentation available in the online help, see [“Solaris Management Console \(Overview\)”](#) in *System Administration Guide: Basic Administration*.

Management Scope

The term *management scope* refers to the name service environment that you choose to use with the selected management tool. The management scope choices for the resource control and performance tools are the `/etc/project` local file, or NIS.

The management scope that you select during a console session should correspond to the primary name service that is identified in the `/etc/nsswitch.conf` file.

Performance Tool

The Performance tool is used to monitor resource utilization. Resource utilization can be summarized for the system, viewed by project, or viewed for an individual user.

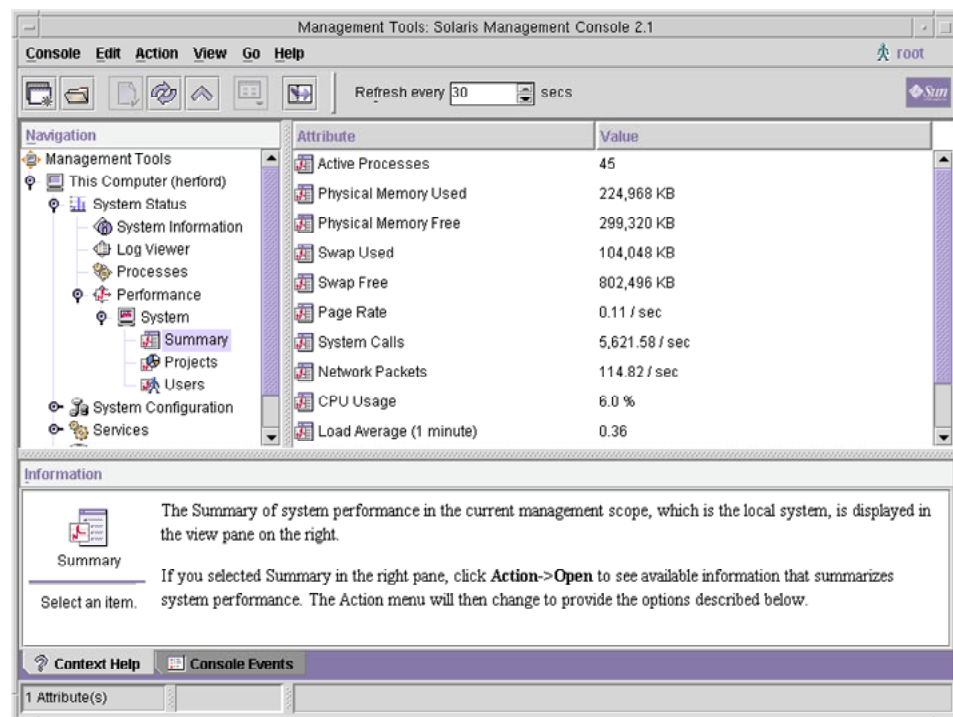


FIGURE 15-1 Performance Tool in the Solaris Management Console

▼ How to Access the Performance Tool

The Performance tool is located under System Status in the Navigation pane. To access the Performance tool, do the following:

1. Click the **System Status** control entity in the **Navigation** pane.

The control entity is used to expand menu items in the Navigation pane.

2. Click the Performance control entity.

3. Click the System control entity.

4. Double-click Summary, Projects, or Users.

Your choice depends on the usage you want to monitor.

Monitoring by System

Values are shown for the following attributes.

Attribute	Description
Active Processes	Number of processes that are active on the system
Physical Memory Used	Amount of system memory that is in use
Physical Memory Free	Amount of system memory that is available
Swap Used	Amount of system swap space that is in use
Swap Free	Amount of free system swap space
Page Rate	Rate of system paging activity
System Calls	Number of system calls per second
Network Packets	Number of network packets that are transmitted per second
CPU Usage	Percentage of CPU that is currently in use
Load Average	Number of processes in the system run queue which are averaged over the last 1, 5, and 15 minutes

Monitoring by Project or User Name

Values are shown for the following attributes.

Attribute	Short Name	Description
Input Blocks	inblk	Number of blocks read

Attribute	Short Name	Description
Blocks Written	oublk	Number of blocks written
Chars Read/Written	ioch	Number of characters read and written
Data Page Fault Sleep Time	dftime	Amount of time spent processing data page faults
Involuntary Context Switches	ictx	Number of involuntary context switches
System Mode Time	stime	Amount of time spent in the kernel mode
Major Page Faults	majfl	Number of major page faults
Messages Received	mrcv	Number of messages received
Messages Sent	msend	Number of messages sent
Minor Page Faults	minf	Number of minor page faults
Num Processes	nprocs	Number of processes owned by the user or the project
Num LWPs	count	Number of lightweight processes
Other Sleep Time	slptime	Sleep time other than tftime, dftime, kftime, and ltime
CPU Time	pctcpu	Percentage of recent CPU time used by the process, the user, or the project
Memory Used	pctmem	Percentage of system memory used by the process, the user, or the project
Heap Size	brksize	Amount of memory allocated for the process data segment
Resident Set Size	rsssize	Current amount of memory claimed by the process
Process Image Size	size	Size of the process image in Kbytes
Signals Received	sig	Number of signals received
Stopped Time	stoptime	Amount of time spent in the stopped state
Swap Operations	swaps	Number of swap operations in progress

Attribute	Short Name	Description
System Calls Made	<code>sysc</code>	Number of system calls made over the last time interval
System Page Fault Sleep Time	<code>kftime</code>	Amount of time spent processing page faults
System Trap Time	<code>ttime</code>	Amount of time spent processing system traps
Text Page Fault Sleep Time	<code>tftime</code>	Amount of time spent processing text page faults
User Lock Wait Sleep Time	<code>ltime</code>	Amount of time spent waiting for user locks
User Mode Time	<code>utime</code>	Amount of time spent in the user mode
User and System Mode Time	<code>time</code>	The cumulative CPU execution time
Voluntary Context Switches	<code>vctx</code>	Number of voluntary context switches
Wait CPU Time	<code>wtime</code>	Amount of time spent waiting for CPU (latency)

Resource Controls Tab

Resource controls allow you to associate a project with a set of resource constraints. These constraints determine the allowable resource usage of tasks and processes that run in the context of the project.

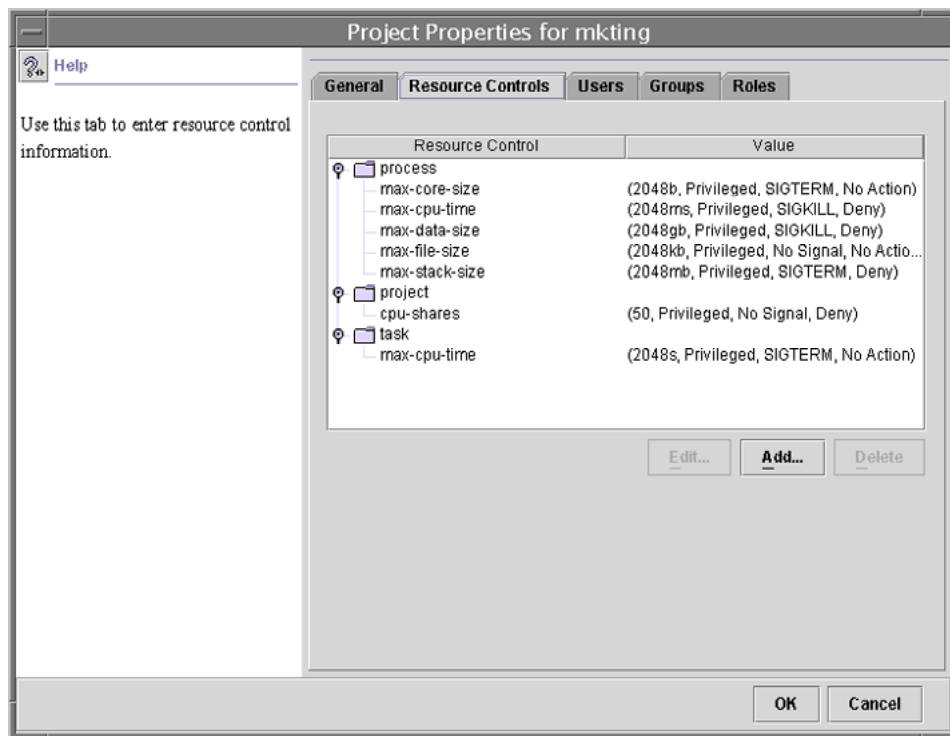


FIGURE 15-2 Resource Controls Tab in the Solaris Management Console

▼ How to Access the Resource Controls Tab

The Resource Controls tab is located under System Configuration in the Navigation pane. To access Resource Controls, do the following:

1. Click the **System Configuration** control entity in the **Navigation** pane.
2. **Double-click** **Projects**.
3. Click on a project in the console main window to select it.
4. Select **Properties** from the **Action** menu.
5. Click the **Resource Controls** tab.

View, add, edit, or delete resource control values for processes, projects, and tasks.

Resource Controls You Can Set

The following table shows the resource controls that can be set in the console. The table describes the resource that is constrained by each control. The table also identifies the default units that are used by the `project` database for that resource. The default units are of two types:

- Quantities represent a limited amount.
- Indexes represent a maximum valid identifier.

Thus, `project.cpu-shares` specifies the number of shares to which the project is entitled. `process.max-file-descriptor` specifies the highest file number that can be assigned to a process by the `open(2)` system call.

TABLE 15-1 Standard Resource Controls Available in the Solaris Management Console

Control Name	Description	Default Unit
<code>project.cpu-shares</code>	The number of CPU shares that are granted to this project for use with the fair share scheduler (FSS) (see the <code>FSS(7)</code> man page)	Quantity (shares)
<code>task.max-cpu-time</code>	Maximum CPU time that is available to this task's processes	Time (seconds)
<code>task.max-lwps</code>	Maximum number of LWPs simultaneously available to this task's processes	Quantity (LWPs)
<code>process.max-cpu-time</code>	Maximum CPU time that is available to this process	Time (seconds)
<code>process.max-file-descriptor</code>	Maximum file descriptor index that is available to this process	Index (maximum file descriptor)
<code>process.max-file-size</code>	Maximum file offset that is available for writing by this process	Size (bytes)
<code>process.max-core-size</code>	Maximum size of a core file that is created by this process	Size (bytes)
<code>process.max-data-size</code>	Maximum heap memory that is available to this process	Size (bytes)
<code>process.max-stack-size</code>	Maximum stack memory segment that is available to this process	Size (bytes)

TABLE 15-1 Standard Resource Controls Available in the Solaris Management Console
(Continued)

Control Name	Description	Default Unit
<code>process.max-address-space</code>	Maximum amount of address space, as summed over segment sizes, available to this process	Size (bytes)

Setting Values

You can view, add, edit, or delete resource control values for processes, projects, and tasks. These operations are performed through dialog boxes in the console.

Resource controls and values are viewed in tables in the console. The Resource Control column lists the resource controls that can be set. The Value column displays the properties that are associated with each resource control. In the table, these values are enclosed in parentheses, and they appear as plain text separated by commas. The values in parentheses comprise an “action clause.” Each action clause is composed of a threshold, a privilege level, one signal, and one local action that is associated with the particular threshold. Each resource control can have multiple action clauses, which are also separated by commas.

Note – On a running system, values that are altered in the `project` database through the console only take effect for new tasks that are started in a project.

Console References

For information on projects and tasks, see [Chapter 2](#). For information on resource controls, see [Chapter 6](#). For information on the fair share scheduler (FSS), see [Chapter 8](#).

Note – Not all resource controls can be set in the console. See [Table 15-1](#) for the list of controls that can be set in the console.

PART II Zones

This part introduces Solaris™ Zones software partitioning technology, which provides a means of virtualizing operating system services to create an isolated environment for running applications. This isolation prevents processes that are running in one zone from monitoring or affecting processes running in other zones.

Introduction to Solaris Zones

The Solaris Zones facility in the Solaris Operating System provides an isolated environment in which to run applications on your system. Solaris Zones are a component of the Solaris Container environment.

This chapter covers the following topics:

- “Zones Overview” on page 195
- “When to Use Zones” on page 196
- “How Zones Work” on page 198
- “Features Provided by Non-Global Zones” on page 203
- “Setting Up Zones on Your System (Task Map)” on page 204

If you are ready to start creating zones on your system, skip to [Chapter 17](#).

Zones Overview

The Solaris Zones partitioning technology is used to virtualize operating system services and provide an isolated and secure environment for running applications. A *zone* is a virtualized operating system environment created within a single instance of the Solaris Operating System. 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 processes that are running in one zone from monitoring or affecting processes that are running in other zones. Even a process running with superuser 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 machine on which they are deployed. Examples of these attributes include physical device paths.

Zones can be used on any machine that is running the Solaris 10 release. The upper limit for the number of zones on a system is 8192. The number of zones that can be effectively hosted on a single system is determined by the total resource requirements of the application software running in all of the zones.

When to Use Zones

Zones are ideal for environments that consolidate a number of applications on a single server. The cost and complexity of managing numerous machines make it advantageous to consolidate several applications on larger, more scalable servers.

The following figure shows a system with four zones. Each of the zones `apps`, `users`, and `work` is running a workload unrelated to the workloads of the other zones, in a sample consolidated environment. This example illustrates that different versions of the same application can be run without negative consequences in different zones, to match the consolidation requirements. Each zone can provide a customized set of services.

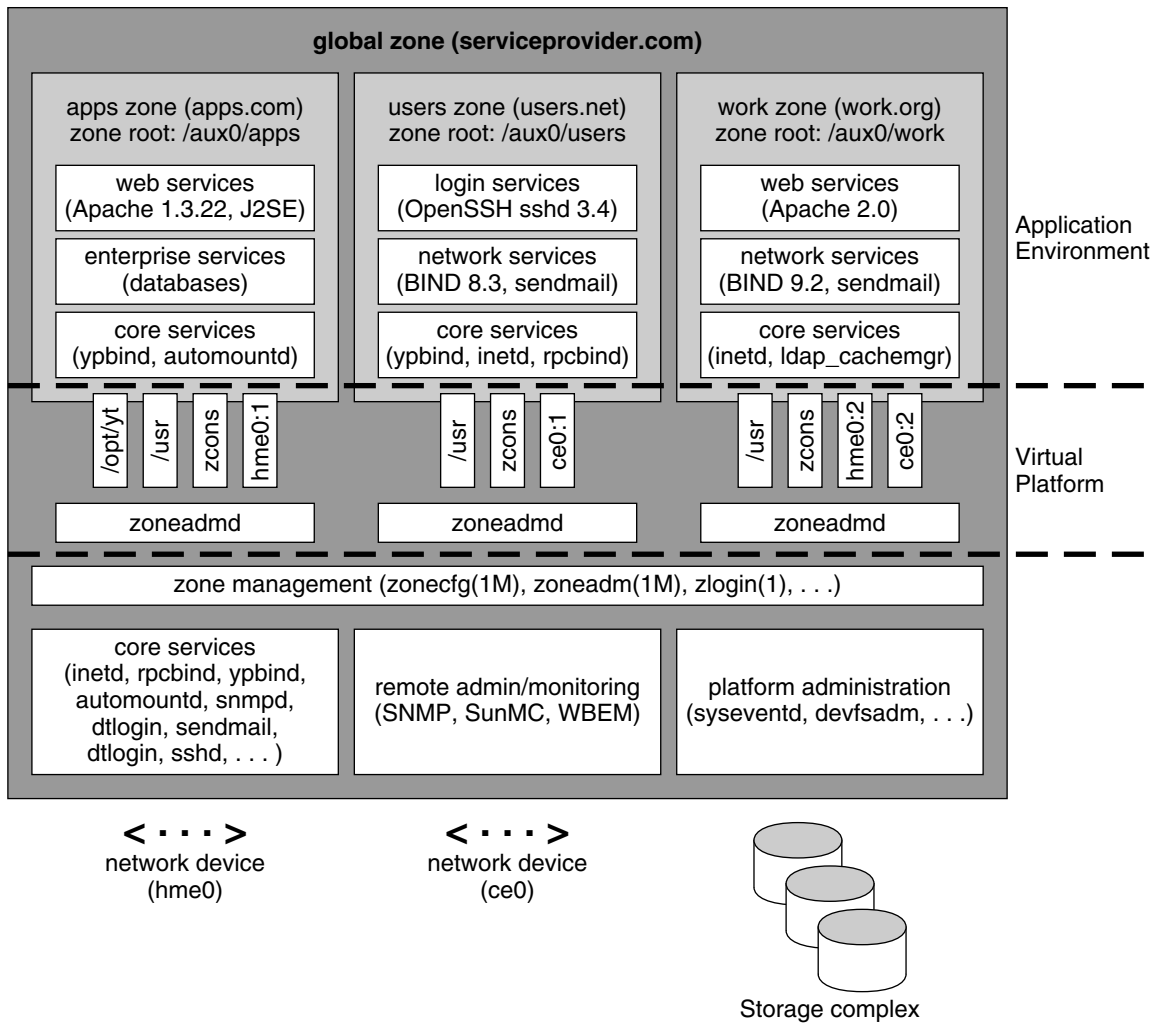


FIGURE 16-1 Zones Server Consolidation Example

Zones enable more efficient resource utilization on your system. Dynamic resource reallocation permits unused resources to be shifted to other containers as needed. Fault and security isolation mean that poorly behaved applications do not require a dedicated and underutilized system. With the use of zones, these applications can be consolidated with other applications.

Zones allow you to delegate some administrative functions while maintaining overall system security.

How Zones Work

A non-global zone can be thought of as a box. One or more applications can run in this box without interacting with the rest of the system. Solaris zones isolate software applications or services by using flexible, software-defined boundaries. Applications that are running in the same instance of the Solaris Operating System can then be managed independently of one other. Thus, different versions of the same application can be run in different zones, to match the requirements of your configuration.

Each zone that requires network connectivity has one or more dedicated IP addresses. A process assigned to a zone can manipulate, monitor, and directly communicate with other processes that are assigned to the same zone. The process cannot perform these functions with processes that are assigned to other zones in the system or with processes that are not assigned to a zone. Processes that are assigned to different zones are only able to communicate through network APIs.

Every Solaris system contains a *global zone*. The global zone has a dual function. The global zone is both the default zone for the system and the zone used for system-wide administrative control. All processes run in the global zone if no *non-global zones*, referred to simply as zones, are created by the *global administrator*.

The global zone is the only zone from which a non-global zone can be configured, installed, managed, or uninstalled. Only the global zone is bootable from the system hardware. Administration of the system infrastructure, such as physical devices, routing, or dynamic reconfiguration (DR), is only possible in the global zone. Appropriately privileged processes running in the global zone can access objects associated with other zones.

Unprivileged processes in the global zone might be able to perform operations not allowed to privileged processes in a non-global zone. For example, users in the global zone can view information about every process in the system. If this capability presents a problem for your site, you can restrict access to the global zone.

Each zone, including the global zone, is assigned a zone name. The global zone always has the name `global`. Each zone is also given a unique numeric identifier, which is assigned by the system when the zone is booted. The global zone is always mapped to ID 0. Zone names and numeric IDs are discussed in [“Using the `zonecfg` Command” on page 209](#).

Each zone also has a node name that is completely independent of the zone name. The node name is assigned by the administrator of the zone. For more information, see [“Non-Global Zone Node Name” on page 295](#).

Each zone has a path to its root directory that is relative to the global zone’s root directory. For more information, see [“Using the `zonecfg` Command” on page 209](#).

The scheduling class for a non-global zone is set to the scheduling class for the system.

You can also set the scheduling class for a zone through the dynamic resource pools facility. If the zone is associated with a pool that has its `pool.scheduler` property set to a valid scheduling class, then processes running in the zone run in that scheduling class by default. To associate a zone with a resource pool, see [Step 6 in “How to Configure the Zone” on page 226](#). For information on assigning a scheduling class to a resource pool, see [“Introduction to Resource Pools” on page 133](#) and [“How to Associate a Pool With a Scheduling Class” on page 163](#).

You can use the `priocntl` described in the `priocntl(1)` man page to move running processes into a different scheduling class without changing the default scheduling class and rebooting.

Summary of Zone Features

The following table summarizes the characteristics of global and non-global zones.

Type of Zone	Characteristic
Global	<ul style="list-style-type: none"> ■ Is assigned ID 0 by the system ■ Provides the single instance of the Solaris kernel that is bootable and running on the system ■ Contains a complete installation of the Solaris system software packages ■ Can contain additional software packages or additional software, directories, files, and other data not installed through packages ■ Provides a complete and consistent product database that contains information about all software components installed in the global zone ■ Holds configuration information specific to the global zone only, such as the global zone host name and file system table ■ Is the only zone that is aware of all devices and all file systems ■ Is the only zone with knowledge of non-global zone existence and configuration ■ Is the only zone from which a non-global zone can be configured, installed, managed, or uninstalled

Type of Zone	Characteristic
Non-Global	<ul style="list-style-type: none"> ■ Is assigned a zone ID by the system when the zone is booted ■ Shares operation under the Solaris kernel booted from the global zone ■ Contains an installed subset of the complete Solaris Operating System software packages ■ Contains Solaris software packages shared from the global zone ■ Can contain additional installed software packages not shared from the global zone ■ Can contain additional software, directories, files, and other data created on the non-global zone that are not installed through packages or shared from the global zone ■ Has a complete and consistent product database that contains information about all software components installed on the zone, whether present on the non-global zone or shared read-only from the global zone ■ Is not aware of the existence of any other zones ■ Cannot install, manage, or uninstall other zones, including itself ■ Has configuration information specific to that non-global zone only, such as the non-global zone host name and file system table

How Non-Global Zones Are Administered

A global administrator has superuser privileges or the Primary Administrator role. When logged in to the global zone, the global administrator can monitor and control the system as a whole.

A non-global zone can be administered by a *zone administrator*. The global administrator assigns the Zone Management profile to the zone administrator. The privileges of a zone administrator are confined to a non-global zone.

How Non-Global Zones Are Created

The global administrator uses the `zonecfg` command to configure a zone by specifying various parameters for the zone's virtual platform and application environment. The zone is then installed by the global administrator, who uses the zone administration command `zoneadm` to install software at the package level into the file system hierarchy established for the zone. The global administrator can log in to the installed zone by using the `zlogin` command. At first login, the internal configuration for the zone is completed. The `zoneadm` command is then used to boot the zone.

For information about zone configuration, see [Chapter 17](#). For information about zone installation, see [Chapter 19](#). For information about zone login, see [Chapter 21](#).

Non-Global Zone State Model

A non-global zone can be in one of the following six states:

Configured	The zone's configuration is complete and committed to stable storage. However, those elements of the zone's application environment that must be specified after initial boot are not yet present.
Incomplete	During an install or uninstall operation, <code>zoneadm</code> sets the state of the target zone to incomplete. Upon successful completion of the operation, the state is set to the correct state.
Installed	The zone's configuration is instantiated on the system. The <code>zoneadm</code> command is used to verify that the configuration can be successfully used on the designated Solaris system. Packages are installed under the zone's root path. In this state, the zone has no associated virtual platform.
Ready	The virtual platform for the zone is established. The kernel creates the <code>zsched</code> process, network interfaces are plumbed, file systems are mounted, and devices are configured. A unique zone ID is assigned by the system. At this stage, no processes associated with the zone have been started.
Running	User processes associated with the zone application environment are running. The zone enters the running state as soon as the first user process associated with the application environment (<code>init</code>) is created.
Shutting down and Down	These states are transitional states that are visible while the zone is being halted. However, a zone that is unable to shut down for any reason will stop in one of these states.

[Chapter 20](#) and the `zoneadm(1M)` man page describe how to use the `zoneadm` command to initiate transitions between these states.

TABLE 16-1 Commands That Affect Zone State

Current Zone State	Applicable Commands
Configured	<code>zonecfg -z <i>zonename</i> verify</code> <code>zonecfg -z <i>zonename</i> commit</code> <code>zonecfg -z <i>zonename</i> delete</code> <code>zoneadm -z <i>zonename</i> verify</code> <code>zoneadm -z <i>zonename</i> install</code>
Incomplete	<code>zoneadm -z <i>zonename</i> uninstall</code>
Installed	<code>zoneadm -z <i>zonename</i> ready</code> (optional) <code>zoneadm -z <i>zonename</i> boot</code> <code>zoneadm -z <i>zonename</i> uninstall</code> uninstalls the configuration of the specified zone from the system.
Ready	<code>zoneadm -z <i>zonename</i> boot</code> <code>zoneadm halt</code> and system reboot return a zone in the ready state to the installed state.
Running	<code>zlogin <i>options</i> <i>zonename</i></code> <code>zoneadm -z <i>zonename</i> reboot</code> <code>zoneadm -z <i>zonename</i> halt</code> returns a ready zone to the installed state. <code>zoneadm halt</code> and system reboot return a zone in the running state to the installed state.

Non-Global Zone Characteristics

A zone provides isolation at almost any level of granularity you require. A zone does not need a dedicated CPU, a physical device, or a portion of physical memory. These resources can either be multiplexed across a number of zones running within a single domain or system, or allocated on a per-zone basis using the resource management features available in the operating system.

Each zone can provide a customized set of services. To enforce basic process isolation, a process can see or signal only those processes that exist in the same zone. Basic communication between zones is accomplished by giving each zone at least one logical network interface. An application running in one zone cannot observe the network traffic of another zone. This isolation is maintained even though the respective streams of packets travel through the same physical interface.

Each zone is given a portion of the file system hierarchy. Because each zone is confined to its subtree of the file system hierarchy, a workload running in a particular zone cannot access the on-disk data of another workload running in a different zone.

Files used by naming services reside within a zone's own root file system view. Thus, naming services in different zones are isolated from one other and the services can be configured differently.

Using Resource Management Features With Non-Global Zones

If you use resource management features, you should align the boundaries of the resource management controls with those of the zones. This alignment creates a more complete model of a virtual machine, where namespace access, security isolation, and resource usage are all controlled.

Any special requirements for using the various resource management features with zones are addressed in the individual chapters of this manual that document those features.

Features Provided by Non-Global Zones

Non-global zones provide the following features:

- | | |
|----------------|---|
| Security | Once a process has been placed in a zone other than the global zone, neither the process nor any of its subsequent children can change zones.

Network services can be run in a zone. By running network services in a zone, you limit the damage possible in the event of a security violation. An intruder who successfully exploits a security flaw in software running within a zone is confined to the restricted set of actions possible within that zone. The privileges available within a zone are a subset of those available in the system as a whole. |
| Isolation | Zones allow the deployment of multiple applications on the same machine, even if those applications operate in different trust domains, require exclusive access to a global resource, or present difficulties with global configurations. For example, multiple applications running in different zones on the same system can bind to the same network port by using the distinct IP addresses associated with each zone or by using the wildcard address. The applications are also prevented from monitoring or intercepting each other's network traffic, file system data, or process activity. |
| Virtualization | Zones provide a virtualized environment that can hide details such as physical devices and the system's primary IP address and host name from applications. The same application environment can be |

maintained on different physical machines. The virtualized environment allows separate administration of each zone. Actions taken by a zone administrator in a non-global zone do not affect the rest of the system.

- Granularity** A zone can provide isolation at almost any level of granularity. See [“Non-Global Zone Characteristics” on page 202](#) for more information.
- Environment** Zones do not change the environment in which applications execute except when necessary to achieve the goals of security and isolation. Zones do not present a new API or ABI to which applications must be ported. Instead, zones provide the standard Solaris interfaces and application environment, with some restrictions. The restrictions primarily affect applications that attempt to perform privileged operations.

Applications in the global zone run without modification, whether or not additional zones are configured.

Setting Up Zones on Your System (Task Map)

The following table provides a basic overview of the tasks that are involved in setting up zones on your system for the first time.

Task	Description	For Instructions
Identify the applications that you would like to run in zones.	Review the applications running on your system: <ul style="list-style-type: none">■ Determine which applications are critical to your business goals.■ Assess the system needs of the applications you are running.	Refer to your business goals and to your system documentation if necessary.

Task	Description	For Instructions
Determine how many zones to configure.	Assess: <ul style="list-style-type: none"> ■ The performance requirements of the applications you intend to run in zones ■ The availability of the recommended 100 MB of free disk space per zone to be installed (Optional) If you are also using resource management features on your system, align the zones with the resource management boundaries	See “Evaluating the Current System Setup” on page 221, Chapter 1.
Perform the preconfiguration tasks.	Determine the zone name and the zone path, obtain IP addresses, and determine the required file systems and devices for each zone. (Optional) To set a default scheduling class for the non-global zone that is different from the system default, you can associate the zone with a pool that specifies a default scheduler. The scheduling class is set at the pool level through the <code>pool.scheduler</code> property.	For information on the zone name and path, IP addresses, file systems, and devices, see Chapter 17 and “Evaluating the Current System Setup” on page 221. For information on resource pool association, see “How Zones Work” on page 198 and “How to Configure the Zone” on page 226.
Develop configurations.	Configure non-global zones.	See “Configuring, Verifying, and Committing a Zone” on page 226 and the <code>zonecfg(1M)</code> man page.
As global administrator, verify and install configured zones.	Zones must be verified and installed prior to login.	See Chapter 19 and Chapter 20 .
As global administrator, log in to the non-global zone.	Log in to perform the initial internal configuration of the zone, including assigning the zone root password.	See Chapter 21 and Chapter 22 .
As global administrator, boot the non-global zone.	Boot the zone to place the zone in the running state.	See Chapter 19 and Chapter 20 .

Task	Description	For Instructions
Prepare the new zone for production use.	Create user accounts, add additional software, and customize the zone's configuration.	Refer to the documentation you use to set up a newly installed machine. Special considerations applicable to the zones environment are covered in this guide.

Non-Global Zone Configuration (Overview)

This chapter provides an introduction to non-global zone configuration.

The following topics are covered in this chapter:

- “Pre-Installation Configuration Process” on page 207
- “Zone Components” on page 208
- “Using the `zonecfg` Command” on page 209
- “`zonecfg` Modes” on page 210
- “Zone Configuration Data” on page 212

After you have learned about zone configuration, go to [Chapter 18](#) to configure non-global zones for installation on your system.

Pre-Installation Configuration Process

Before you can install a non-global zone and use it on your system, the zone must be configured.

The `zonecfg` command is used to create the configuration and to determine whether the specified resources and properties are valid on a hypothetical system. The check performed by `zonecfg` for a given configuration verifies the following:

- Ensures that a zone path is specified
- Ensures that all of the required properties for each resource are specified

For more information about the `zonecfg` command, see the `zonecfg(1M)` man page.

Zone Components

This section covers the zone resources and properties that can be configured.

Zone Name and Path

You must choose a name and a path for your zone.

Zone Interfaces

Each zone that requires network connectivity must have one or more dedicated IP addresses. These addresses are associated with logical network interfaces. Zone interfaces configured by the `zonecfg` command will automatically be plumbed and placed in the zone when it is booted.

The `ifconfig` command can be used from the global zone to add or remove logical interfaces in a running zone. For more information, see [“Network Interfaces” on page 301](#).

File Systems Mounted in Zones

Generally, the file systems mounted in a zone include the following:

- The set of file systems mounted when the virtual platform is initialized
- The set of file systems mounted from within the application environment itself

This can include, for example, the following file systems:

- File systems specified in a zone's `/etc/vfstab` file
- AutoFS and AutoFS-triggered mounts
- Mounts explicitly performed by a zone administrator

Certain restrictions are placed on mounts performed from within the application environment. These restrictions prevent the zone administrator from denying service to the rest of the system, or otherwise negatively impacting other zones.

There are security restrictions associated with mounting certain file systems from within a zone. Other file systems exhibit special behavior when mounted in a zone. See [“File Systems and Non-Global Zones” on page 295](#) for more information.

Configured Devices in Zones

The `zonectfg` command uses a rule-matching system to specify which devices should appear in a particular zone. Devices matching one of the rules are included in the zone's `/dev` file system. For more information, see [“How to Configure the Zone” on page 226](#).

Zone-Wide Resource Controls

The global administrator can set privileged zone-wide resource controls for a zone. Zone-wide resource controls limit the total resource usage of all process entities within a zone, regardless of project. These limits are specified in the `zonectfg` configuration. For more information, see [“How to Configure the Zone” on page 226](#).

Including a Comment for a Zone

You can add a comment for a zone by using the `attr` resource type. For more information, see [“How to Configure the Zone” on page 226](#).

Using the `zonectfg` Command

The `zonectfg` command, which is described in the `zonectfg(1M)` man page, is used to configure a zone. The `zonectfg` command can be used in interactive mode, in command-line mode, or in command-file mode. The following operations can be performed using this command:

- Create or delete (destroy) a zone configuration
- Add resources to a particular configuration
- Set properties for resources added to a configuration
- Remove resources from a particular configuration
- Query or verify a configuration
- Commit to a configuration
- Revert to a previous configuration
- Exit from a `zonectfg` session

The `zonectfg` prompt is of the following form:

```
zonectfg:zonename>
```

When you are configuring a specific resource type, such as a file system, that resource type is also included in the prompt:

```
zoncfg:zonename:fs>
```

For more information, including procedures that show how to use the various `zoncfg` components described in this chapter, see [Chapter 18](#).

zoncfg Modes

The concept of a *scope* is used for the user interface. The scope can be either *global* or *resource specific*. The default scope is *global*.

In the *global* scope, the `add` subcommand and the `select` subcommand are used to select a specific resource. The scope then changes to that resource type.

- For the `add` subcommand, the `end` or `cancel` subcommands are used to complete the resource specification.
- For the `select` subcommand, the `end` or `cancel` subcommands are used to complete the resource modification.

The scope then reverts back to *global*.

Certain subcommands, such as `add`, `remove`, and `set`, have different semantics in each scope.

zoncfg Interactive Mode

In interactive mode, the following subcommands are supported. For detailed information about semantics and options used with the subcommands, see the `zoncfg(1M)` man page for options. For any subcommand that could result in destructive actions or loss of work, the system requests user confirmation before proceeding. You can use the `-F` (*force*) option to bypass this confirmation.

`help` Print general help, or display help about a given resource.

```
zoncfg:my-zone:inherit-pkg-dir> help
```

`create` Begin configuring an in-memory configuration for the specified new zone for one of these purposes:

- To apply the Sun default settings to a new configuration. This method is the default.
- With the `-t template` option, to create a configuration that is identical to the specified template. The zone name is changed from the template name to the new zone name.
- With the `-F` option, to overwrite an existing configuration.

	<ul style="list-style-type: none"> ■ With the <code>-b</code> option, to create a blank configuration in which nothing is set.
<code>export</code>	Print the configuration to <code>stdout</code> , or to the output file specified, in a form that can be used in a command file.
<code>add</code>	<p>In the global scope, add the specified resource type to the configuration.</p> <p>In the resource scope, add a property of the given name with the given value.</p> <p>See “How to Configure the Zone” on page 226 and the <code>zonectfg(1M)</code> man page for more information.</p>
<code>set</code>	Set a given property name to the given property value. Note that some properties, such as <code>zonepath</code> , are global, while others are resource specific. Thus, this command is applicable in both the global and resource scopes.
<code>select</code>	Applicable only in the global scope. Select the resource of the given type that matches the given property name-property value pair criteria for modification. The scope is changed to that resource type. You must specify a sufficient number of property name-value pairs for the resource to be uniquely identified.
<code>remove</code>	<p>In the global scope, remove the specified resource type. You must specify a sufficient number of property name-value pairs for the resource type to be uniquely identified.</p> <p>In the resource scope, remove the specified property name-property value from the current resource.</p>
<code>end</code>	<p>Applicable only in the resource scope. End the resource specification.</p> <p>The <code>zonectfg</code> command then verifies that the current resource is fully specified.</p> <ul style="list-style-type: none"> ■ If the resource is fully specified, it is added to the in-memory configuration and the scope will revert back to global. ■ If the specification is incomplete, the system displays an error message that describes what needs to be.
<code>cancel</code>	Applicable only in the resource scope. End the resource specification and reset the scope to global. Any partially specified resources are not retained.
<code>delete</code>	Destroy the specified configuration. Delete the configuration both from memory and from stable storage. You must use the <code>-F</code> (force) option with <code>delete</code> .



Caution – This action is instantaneous. No commit is required, and a deleted zone cannot be reverted.

- `info` Display information about the current configuration or the global resource properties `zonpath`, `autoboot`, and `pool`. If a resource type is specified, display information only about resources of that type. In the resource scope, this subcommand applies only to the resource being added or modified.
- `verify` Verify current configuration for correctness. Ensure that all resources have all of their required properties specified.
- `commit` Commit current configuration from memory to stable storage. Until the in-memory configuration is committed, changes can be removed with the `revert` subcommand. A configuration must be committed to be used by `zoneadm`. This operation is attempted automatically when you complete a `zonecfg` session. Because only a correct configuration can be committed, the commit operation automatically does a `verify`.
- `revert` Revert configuration back to the last committed state.
- `exit -F` Exit the `zonecfg` session. You can use the `-F` (force) option with `exit`.

A `commit` is automatically attempted if needed. Note that an EOF character can also be used to exit the session.

zonecfg Command-File Mode

In command-file mode, input is taken from a file. The `export` subcommand described in “[zonecfg Interactive Mode](#)” on page 210 is used to produce this file. The configuration can be printed to standard output, or the `-f` option can be used to specify an output file.

Zone Configuration Data

Zone configuration data consists of two kinds of entities: resources and properties. Each resource has a type, and each resource can also have a set of one or more properties. The properties have names and values. The set of properties is dependent on the resource type.

Resource Types

The resource types are described as follows:

Zone name	<p>The zone name identifies the zone to the configuration utility. The following rules apply to zone names:</p> <ul style="list-style-type: none">■ Each zone must have a unique name.■ A zone name is case-sensitive.■ A zone name must begin with an alpha-numeric character. <p>The name can contain alpha-numeric characters, underbars (<code>_</code>), hyphens (<code>-</code>), and periods (<code>.</code>).</p> <ul style="list-style-type: none">■ The name cannot be longer than 64 characters.■ The name <code>global</code> and all names beginning with <code>SUNW</code> are reserved and cannot be used.
zonepath	<p>The zone path resource is the path to the zone root. Each zone has a path to its root directory that is relative to the global zone's root directory. At installation time, the global zone directory is required to have restricted visibility. It must be owned by <code>root</code> with the mode <code>700</code>.</p> <p>The non-global zone's root path is one level lower. The zone's root directory has the same ownership and permissions as the root directory in the global zone. The zone directory must be owned by <code>root</code> with the mode <code>755</code>. This hierarchy ensures that unprivileged users in the global zone are prevented from traversing a non-global zone's file system. See "Traversing File Systems" on page 300 for a further discussion of this issue.</p>
fs	<p>Each zone can have various file systems that are mounted when the zone transitions from the installed state to the ready state. The file system resource specifies the path to the file system mount point. For more information about the use of file systems in zones, see "File Systems and Non-Global Zones" on page 295.</p>
inherit-pkg-dir	<p>The <code>inherit-pkg-dir</code> resource is used to represent directories that contain packaged software that a non-global zone shares with the global zone.</p>

The contents of software packages transferred into the `inherit-pkg-dir` directory are inherited in read-only mode by the non-global zone. The zone's packaging database is updated to reflect the packages. These resources cannot be modified or removed after the zone has been installed using `zoneadm`.

Note – Four default `inherit-pkg-dir` resources are included in the configuration. These directory resources indicate which directories should have their associated packages inherited from the global zone. The resources are implemented through a read-only loopback file system mount.

- `/lib`
 - `/platform`
 - `/sbin`
 - `/usr`
-

<code>net</code>	The network interface resource is the virtual interface name. Each zone can have network interfaces that should be plumbed when the zone transitions from the installed state to the ready state.
<code>device</code>	The device resource is the device matching specifier. Each zone can have devices that should be configured when the zone transitions from the installed state to the ready state.
<code>rctl</code>	The <code>rctl</code> resource is used for zone-wide resource controls. The controls should be enabled when the zone transitions from the installed state to the ready state. The zone-wide resource controls implemented in this release are <code>zone.cpu-shares</code> and <code>zone.max-lwps</code> .
<code>attr</code>	This generic attribute can be used for user comments or by other subsystems. The name property of an <code>attr</code> must begin with an alpha-numeric character. The name property can contain alpha-numeric characters, hyphens (-), and periods (.). Attribute names beginning with <code>zone.</code> are reserved for use by the system.

Resource Type Properties

Some resource types also have properties to configure. The following properties are associated with the resource types shown.

fs dir, special, raw, type, options

The lines in the following example specify that /dev/dsk/c0t0d0s2 in the global zone is to be mounted as /mnt in a zone being configured. The raw property specifies an optional device on which the fsck command is to be run before an attempt is made to mount the file system. The file system type to use is UFS. The options nodevices and logging are added.

```
zonecfg:my-zone> add fs
zonecfg:my-zone:fs> set dir=/mnt
zonecfg:my-zone:fs> set special=/dev/dsk/c0t0d0s2
zonecfg:my-zone:fs> set raw=/dev/rdisk/c0t0d0s2
zonecfg:my-zone:fs> set type=ufs
zonecfg:my-zone:fs> add options [nodevices,logging]
zonecfg:my-zone:fs> end
```

For more information, see [“The -o nosuid Option” on page 295](#), [“Security Restrictions and File System Behavior” on page 298](#), and the fsck(1M) and mount(1M) man pages. Also note that section 1M man pages are available for mount options that are unique to a specific filesystem. The names of these man pages have the form mount_ filesystem.

inherit-pkg-dir dir

The lines in the following example specify that /opt/sfw is to be loopback mounted from the global zone.

```
zonecfg:my-zone> add inherit-pkg-dir
zonecfg:my-zone:inherit-pkg-dir> set dir=/opt/sfw
zonecfg:my-zone:inherit-pkg-dir> end
```

net address, physical

In the following example, IP address 192.168.0.1 is added to a zone. An hme0 card is used for the physical interface.

```
zonecfg:my-zone> add net
zonecfg:my-zone:net> set physical=hme0
zonecfg:my-zone:net> set address=192.168.0.1
zonecfg:my-zone:net> end
```

Note – To determine which physical interface to use, type `ifconfig -a` on your system. Each line of the output, other than loopback driver lines, begins with the name of a card installed on your system. Lines that contain LOOPBACK in the descriptions do not apply to cards.

device

match

In the following example, a `/dev/pts` device is included in a zone.

```
zonecfg:my-zone> add device
zonecfg:my-zone:device> set match=/dev/pts*
zonecfg:my-zone:device> end
```

rctl

name, value

In this release, there are two zone-wide resource controls, `zone.cpu-shares` and `zone.max-lwps`.

The `zone.cpu-shares` resource control sets a limit on the number of fair share scheduler (FSS) CPU shares for a zone. 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. For more information, see [“Using the Fair Share Scheduler in a Zones Environment”](#) on page 323.

The `zone.max-lwps` resource control enhances resource isolation by preventing 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.

Note that zone-wide resource control entries in a zone are configured differently than resource control entries in the `project` database. In a zone configuration, the `rctl` resource type consists of three name/value pairs. The names are `priv`, `limit`, and `action`. Each of the names takes a simple value.

```
zonecfg:my-zone> add rctl
zonecfg:my-zone:rctl> set name=zone.max-lwps
zonecfg:my-zone:rctl> add value (priv=privileged,limit=100,action=deny)
zonecfg:my-zone:rctl> end
```

For general information about resource controls and attributes, see [Chapter 6](#) and [“Resource Controls Used in Non-Global Zones”](#) on page 305.

attr

name, type, value

In the following example, a comment about a zone is added.

```
zonecfg:my-zone> add attr
zonecfg:my-zone:attr> set name=comment
zonecfg:my-zone:attr> set type=string
zonecfg:my-zone:attr> set value="Production zone"
```



```
zonecfg:my-zone:attr> end
```

You can use the `export` subcommand to print a zone configuration to `stdout`. The configuration is saved in a form that can be used in a command file.

Tecla Command-Line Editing Library

The Tecla command-line editing library is included for use with the `zonecfg` command. The library provides a mechanism for command-line history and editing support.

The Tecla command-line editing library is documented in the following man pages:

- `enhance(1)`
- `libtecla(3LIB)`
- `ef_expand_file(3TECLA)`
- `gl_get_line(3TECLA)`
- `gl_io_mode(3TECLA)`
- `pca_lookup_file(3TECLA)`
- `tecla(5)`

Planning and Configuring Non-Global Zones (Tasks)

This chapter describes what you need to do before you can configure a zone on your system. This chapter also describes how to configure a zone, modify a zone configuration, and delete a zone configuration from your system.

For an introduction to the zone configuration process, see [Chapter 17](#).

Planning and Configuring a Non-Global Zone (Task Map)

Before you set up your system to use zones, you must first collect information and make decisions about how to configure the zones. The following task map summarizes how to plan and configure a zone.

Task	Description	For Instructions
Plan your zone strategy.	<ul style="list-style-type: none"> ■ Evaluate the applications running on your system to determine which applications you want to run in a zone. ■ Assess the availability of disk space to hold the files that are unique in the zone. ■ If you are also using resource management features, determine how to align the zone with the resource management boundaries. 	Refer to historical usage. Also see “Disk Space Requirements” on page 221 and “Resource Pools Used in Zones” on page 135.
Determine the name for the zone.	Decide what to call the zone based on the naming conventions.	See “Zone Configuration Data” on page 212 and “Zone Host Name” on page 223.
Obtain or configure IP addresses for the zone.	Depending on your configuration, you must obtain at least one IP address for each non-global zone that you want to have network access.	See “Determine the Zone Host Name and Obtain the Network Address” on page 223 and <i>System Administration Guide: IP Services</i> .
Determine which file systems you want to mount in the zone.	Review your application requirements.	See “File Systems Mounted in Zones” on page 208 for more information.
Determine which network interfaces should be plumbed in the zone.	Review your application requirements.	See “Network Interfaces” on page 301 for more information.
Determine which devices should be configured in each zone.	Review your application requirements.	Application chapter: ref to come.
Determine the zone path.	Each zone has a path to its root directory that is relative to the global zone’s root directory.	See “Zone Configuration Data” on page 212.
Configure the zone.	Use <code>zonectfg</code> to create a configuration for the zone.	See “Configuring, Verifying, and Committing a Zone” on page 226.

Task	Description	For Instructions
Verify and commit the configured zone.	Determine whether the resources and properties specified are valid on a hypothetical system.	See “Configuring, Verifying, and Committing a Zone” on page 226.

Evaluating the Current System Setup

Zones can be used on any machine that runs the Solaris 10 release. The following primary machine considerations are associated with the use of zones.

- The performance requirements of the applications running within each zone.
- The availability of disk space to hold the files that are unique within each zone.

Disk Space Requirements

There are no limits on how much disk space can be consumed by a zone. The global administrator is responsible for space restriction. The global administrator must ensure that local storage is sufficient to hold a non-global zone’s root file system. Even a small uniprocessor system can support a number of zones running simultaneously.

The nature of the packages installed in the global zone affects the space requirements of the non-global zones that are created. The number of packages and space requirements are factors.

There are two types of non-global zone root file system models: sparse and whole root.

Sparse Root Zones

The sparse root zone model optimizes the sharing of objects in the following ways:

- Only a subset of the root packages, those that have the `pkginfo` parameter `SUNW_PKGTYPE` set to `root`, are installed. See the `pkginfo(4)` man page.
- Read-only loopback file systems are used to gain access to other files.

In this model, only certain root packages are installed in the non-global zone. This includes a subset of the required root packages that are normally installed in the global zone, and any additional root packages that the global administrator might select. Access to other files will be through read-only loopback file systems, identified as `inherit-pkg-dir` resources.

Thus, zones that have `inherit-pkg-dir` resources are called sparse root zones.

- As a general guideline, a zone requires about 100 megabytes of free disk space per zone when the global zone has been installed with all of the standard Solaris packages.
- By default, any additional packages installed in the global zone also populate the non-global zones. The amount of disk space required might be increased accordingly, depending on whether the additional packages deliver files that reside in the `inherit-pkg-dir` resource space.

An additional 40 megabytes of RAM per zone are suggested, but not required on a machine with sufficient swap space.

Whole Root Zones

The whole root model provides the maximum configurability. All of the required and any selected optional Solaris packages are installed into the private file systems of the zone. The advantages of this model include the capability for global zone administrators to customize their zones file system layout. This would be done, for example, to add arbitrary unbundled or third-party packages.

The disk requirements for this model are determined by the disk space used by the packages currently installed in the global zone.

Note – If you create a sparse root zone that contains the following `inherit-pkg-dir` directories:

- `/lib`
- `/platform`
- `/sbin`
- `/usr`

You must remove these directories from the non-global zone's configuration before the zone is installed to have a whole root zone. See [“How to Configure the Zone”](#) on page 226.

Restricting Zone Size

The following options can be used to restrict zone size:

- You can place the zone on a `lofi`-mounted partition. This action will limit the amount of space consumed by the zone to that of the file used by `lofi`. For more information, see the `lofiadm(1M)` and `lofi(7D)` man pages.
- You can use soft partitions to divide disk slices or logical volumes into partitions. You can use these partitions as zone roots, and thus limit per-zone disk consumption. The soft partition limit is 8192 partitions. For more information, see *“Soft Partitions (Overview)”* in *Solaris Volume Manager Administration Guide*.

- You can use the standard partitions of a disk for zone roots, and thus limit per-zone disk consumption.

Determine the Zone Host Name and Obtain the Network Address

You must determine the host name for the zone. Then, you must assign an IPv4 address or manually configure and assign an IPv6 address for the zone if you want it to have network connectivity.

Zone Host Name

The host name you select for the zone must be defined either in the `hosts` database or in the `/etc/inet/ipnodes` database, as specified by the `/etc/nsswitch.conf` file in the global zone. The network databases are files that provide network configuration information. The `nsswitch.conf` file specifies which naming service to use.

If you use local files for the naming service, the `hosts` database is maintained in the `/etc/inet/hosts` file. The host names for zone interfaces are resolved from the local `hosts` database in `/etc/inet/hosts`. Alternatively, the IP address itself can be specified directly when configuring a zone so that no host name resolution is required.

For more information, see “TCP/IP Configuration Files” in *System Administration Guide: IP Services* and “Network Databases and File” in *System Administration Guide: IP Services*.

Zone Network Address

Each zone that requires network connectivity has one or more unique IP addresses. Both IPv4 and IPv6 addresses are supported.

IPv4 Zone Network Address

If you are using IPv4, obtain an address and assign the address to the zone.

A prefix length can also be specified with the IP address. The format of this prefix is *address/prefix-length*, for example, `192.168.1.1/24`. Thus, the address to use is `192.168.1.1` and the netmask to use is `255.255.255.0`, or the mask where the first 24 bits are 1-bits.

IPv6 Zone Network Address

If you are using IPv6, you must manually configure the address. Typically, at least the following two types of addresses must be configured:

Link-local address

A link-local address is of the form `fe80::64-bit interface ID/10`. The `/10` indicates a prefix length of 10 bits.

Address formed from a global prefix configured on the subnet

A global unicast address is based off a 64-bit prefix that the administrator configures for each subnet, and a 64-bit interface ID. The prefix can also be obtained by running the `ifconfig` command with the `-a6` option on any system on the same subnet that has been configured to use IPv6.

The 64-bit interface ID is typically derived from a system's MAC address. For zones use, an alternate address that is unique can be derived from the global zone's IPv4 address as follows:

```
16 bits of zero:upper 16 bits of IPv4 address:lower 16 bits of
IPv4 address:a zone-unique number
```

For example, if the global zone's IPv4 address is 192.168.200.10, a suitable link-local address for a local zone using a zone-unique number of 1 is

```
fe80::c0a8:c80a:1/10. If the global prefix in use on that subnet is
2001:0db8:aabb:ccdd/64, a unique global unicast address for the same
non-global zone is 2001:0db8:aabb:ccdd::c0a8:c80a:1/64. Note that you
must specify a prefix length when configuring an IPv6 address.
```

For more information about link-local and global unicast addresses, see the `inet6(7P)` manual page.

File System Configuration

You can specify a number of mounts to be performed when the virtual platform is set up. File systems that are loopback-mounted into a zone by using the loopback file system (LOFS) virtual file system should be mounted with the `nodevices` option. For information on the `nodevices` option, see [“File Systems and Non-Global Zones” on page 295](#).

LOFS lets you create a new virtual file system so that you can access files by using an alternative path name. In a non-global zone, a loopback mount makes the file system hierarchy look as though it is duplicated under the zone's root. In the zone, all files will be accessible with a path name that starts from the zone's root. LOFS mounting preserves the file system name space.

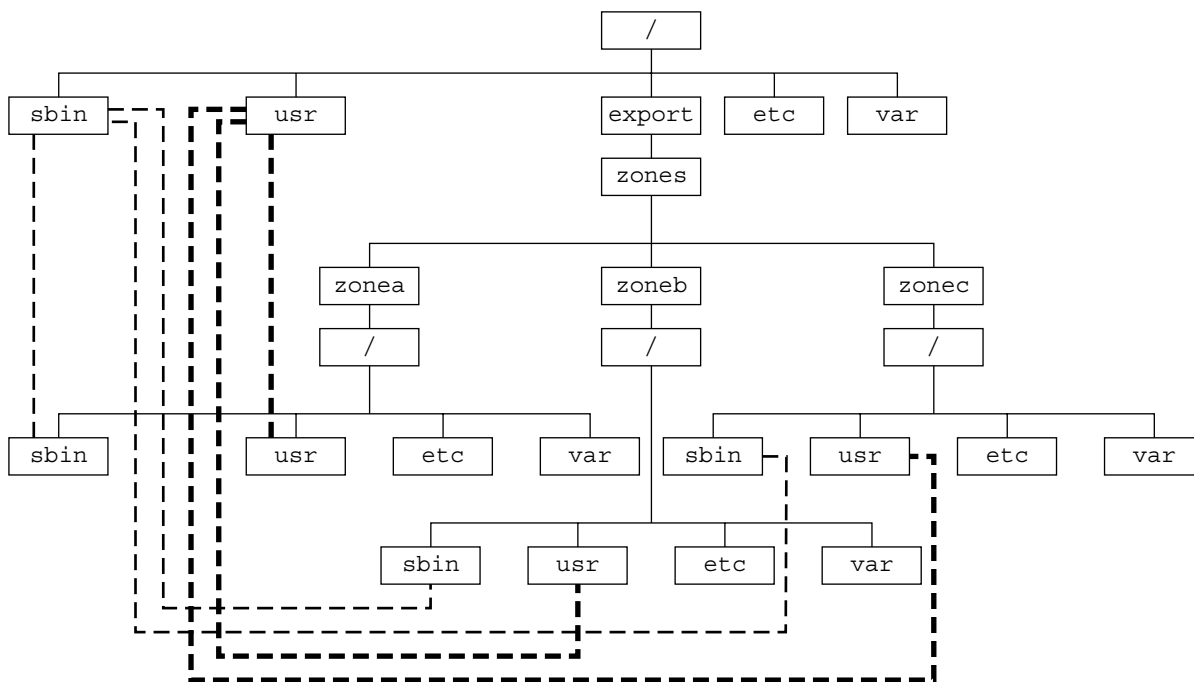


FIGURE 18-1 Loopback-Mounted File Systems

See the `lofs(7S)` man page for more information.

Creating, Revising, and Deleting Non-Global Zone Configurations (Task Map)

Task	Description	For Instructions
Configure a non-global zone.	Use the <code>zonecfg</code> command to create a zone, verify the configuration, and commit the configuration. You can also use a script to configure and boot multiple zones on your system.	“Configuring, Verifying, and Committing a Zone” on page 226, “Script to Configure Multiple Zones” on page 230

Task	Description	For Instructions
Modify a zone configuration.	Use this procedure to modify a resource type in a zone configuration or add a dedicated device to a zone.	“Using the zonecfg Command to Modify a Zone Configuration” on page 232
Revert a zone configuration or delete a zone configuration.	Use the zonecfg command to undo a resource setting made to a zone configuration or to delete a zone configuration.	“Using the zonecfg Command to Revert or Remove a Zone Configuration” on page 234
Delete a zone configuration.	Use the zonecfg command with the delete subcommand to delete a zone configuration from the system.	“How to Delete a Zone Configuration” on page 236

Configuring, Verifying, and Committing a Zone

You use the zonecfg command described in the zonecfg(1M) man page to perform the following actions:

- Create the zone configuration
- Verify that all required information is present
- Commit the non-global zone configuration

While configuring a zone with the zonecfg utility, you can use the revert subcommand to undo the setting for a resource. See [“How to Revert a Zone Configuration” on page 234](#).

A script to configure multiple zones on your system is provided in [“Script to Configure Multiple Zones” on page 230](#).

▼ How to Configure the Zone

You must be the global administrator in the global zone to perform this procedure.

1. Become superuser, or assume the Primary Administrator role.

To create the role and assign the role to a user, see [“Using the Solaris Management Tools With RBAC \(Task Map\)” in *System Administration Guide: Basic Administration*](#).

2. Set up a zone configuration with the zone name you have chosen.

The name `my-zone` is used in this example procedure.

```
global# zonecfg -z my-zone
```

If this is the first time you have configured this zone, you will see the following system message:

```
my-zone: No such zone configured
Use 'create' to begin configuring a new zone.
```

3. Create the new zone configuration.

This procedure uses the Sun default settings.

```
zonecfg:my-zone> create
```

4. Set the zone path, `/export/home/my-zone` in this procedure.

```
zonecfg:my-zone> set zonepath=/export/home/my-zone
```

5. Set the autoboot value.

If set to `true`, the zone is automatically booted when the global zone is booted. The default value is `false`.

```
zonecfg:my-zone> set autoboot=true
```

6. If resource pools are enabled on your system, associate a pool with the zone.

This example uses the default pool, named `pool_default`.

```
zonecfg:my-zone> set pool=pool_default
```

Because a resource pool can have an optional scheduling class assignment, you can use the pools facility to set a default scheduler other than the system default for a non-global zone. For instructions, see [“How to Associate a Pool With a Scheduling Class” on page 163](#) and [“Creating the Configuration” on page 176](#).

7. Add a file system.

```
zonecfg:my-zone> add fs
```

a. Set the mount point for the file system, `/usr/local` in this procedure.

```
zonecfg:my-zone:fs> set dir=/usr/local
```

b. Specify that `/opt/local` in the global zone is to be mounted as `/usr/local` in the zone being configured.

```
zonecfg:my-zone:fs> set special=/opt/local
```

In the non-global zone, the `/usr/local` file system will be readable and writable.

c. Specify the file system type, `lofs` in this procedure.

```
zonecfg:my-zone:fs> set type=lofs
```

d. End the file system specification.

```
zonecfg:my-zone:fs> end
```

This step can be performed more than once to add more than one file system.

8. (Sparse Root Zone Only) Add a shared file system that is loopback-mounted from the global zone.

Do *not* perform this step to create a whole root zone, which does not have any shared file systems. See “Whole Root Zones” on page 222.

```
zonecfg:my-zone> add inherit-pkg-dir
```

a. Specify that /opt/sfw in the global zone is to be mounted in read-only mode in the zone being configured.

```
zonecfg:my-zone:inherit-pkg-dir> set dir=/opt/sfw
```

Note – The zone’s packaging database is updated to reflect the packages. These resources cannot be modified or removed after the zone has been installed using zoneadm.

b. End the inherit-pkg-dir specification.

```
zonecfg:my-zone:inherit-pkg-dir> end
```

This step can be performed more than once to add more than one shared file system.

Note – If you want to create a whole root zone but default shared file systems resources have been added by using inherit-pkg-dir, you must remove these default inherit-pkg-dir resources using zonecfg before you install the zone:

- zonecfg:my-zone> remove inherit-pkg-dir dir=/lib
 - zonecfg:my-zone> remove inherit-pkg-dir dir=/platform
 - zonecfg:my-zone> remove inherit-pkg-dir dir=/sbin
 - zonecfg:my-zone> remove inherit-pkg-dir dir=/usr
-

9. Add a network virtual interface.

```
zonecfg:my-zone> add net
```

a. Set the IP address for the network interface, 192.168.0.1 in this procedure.

```
zonecfg:my-zone:net> set address=192.168.0.1
```

b. Set the physical device type for the network interface, the hme device in this procedure.

```
zonecfg:my-zone:net> set physical=hme0
```

c. End the specification.

```
zonecfg:my-zone:net> end
```

This step can be performed more than once to add more than one network interface.

10. Add a device.

```
zonecfg:my-zone> add device
```

a. Set the device match, /dev/sound/* in this procedure.

```
zonecfg:my-zone:device> set match=/dev/sound/*
```

b. End the device specification.

```
zonecfg:my-zone:device> end
```

This step can be performed more than once to add more than one device.

11. Add a zone-wide resource control.

```
zonecfg:my-zone> add rctl
```

a. Set the name of the resource control, zone.cpu-shares in this procedure.

```
zonecfg:my-zone:rctl> set name=zone.cpu-shares
```

b. Add values for the privilege, the share limit, and the action to be taken when that threshold is reached.

```
zonecfg:my-zone:rctl> add value (priv=privileged,limit=20,action=none)
```

c. End the rctl specification.

```
zonecfg:my-zone:rctl> end
```

This step can be performed more than once to add more than one resource control.

12. Add a comment by using the attr resource type.

```
zonecfg:my-zone> add attr
```

a. Set the name to comment.

```
zonecfg:my-zone:attr> set name=comment
```

b. Set the type to string.

```
zonecfg:my-zone:attr> set type=string
```

c. Set the value to a comment that describes the zone.

```
zonecfg:my-zone:attr> set value="This is my work zone."
```

d. End the attr resource type specification.

```
zonecfg:my-zone:attr> end
```

13. Verify the zone configuration for the zone.

```
zonecfg:my-zone> verify
```

14. Commit the zone configuration for the zone.

```
zonecfg:my-zone> commit
```

15. Exit the zonecfg command.

```
zonecfg:my-zone> exit
```

Note that even if you did not explicitly type `commit` at the prompt, a `commit` is automatically attempted when you type `exit` or an EOF occurs.

See “Installing and Booting Zones” on page 244 to install your committed zone configuration.

Note – The `zonecfg` command also supports multiple subcommands, quoted and separated by semicolons, from the same shell invocation.

```
global# zonecfg -z my-zone "create ; set zonepath=/export/home/my-zone"
```

Script to Configure Multiple Zones

You can use this script to configure and boot multiple zones on your system. The script takes the following parameters:

- The number of zones to be created
- The *zonename* prefix
- The directory to use as the base directory

You must be the global administrator in the global zone to execute the script. The global administrator has superuser privileges in the global zone or assumes the Primary Administrator role.

```
#!/bin/ksh
#
# Copyright 2004 Sun Microsystems, Inc. All rights reserved.
# Use is subject to license terms.
#
#ident      "%Z%M%      %I%      %E% SMI"

if [[ -z "$1" || -z "$2" || -z "$3" ]]; then
    echo "usage: $0 <#-of-zones> <zonename-prefix> <basedir>"
    exit 2
fi

if [[ ! -d $3 ]]; then
    echo "$3 is not a directory"
    exit 1
fi
```

```

fi

nprocs=`psrinfo | wc -l`
nzones=$1
prefix=$2
dir=$3

ip_addrs_per_if=`ndd /dev/ip ip_addrs_per_if`
if [ $ip_addrs_per_if -lt $nzones ]; then
    echo "ndd parameter ip_addrs_per_if is too low ($ip_addrs_per_if)"
    echo "set it higher with 'ndd -set /dev/ip ip_addrs_per_if <num>'"
    exit 1
fi

i=1
while [ $i -le $nzones ]; do
    zoneadm -z $prefix$i list > /dev/null 2>&1
    if [ $? != 0 ]; then
        echo configuring $prefix$i
        F=$dir/$prefix$i.config
        rm -f $F
        echo "create" > $F
        echo "set zonepath=$dir/$prefix$i" >> $F
        zonecfg -z $prefix$i -f $dir/$prefix$i.config 2>&1 | \
            sed 's/^/ /g'
    else
        echo "skipping $prefix$i, already configured"
    fi
    i=`expr $i + 1`
done

i=1
while [ $i -le $nzones ]; do
    j=1
    while [ $j -le $nprocs ]; do
        if [ $i -le $nzones ]; then
            if [ `zoneadm -z $prefix$i list -p | \
                cut -d':' -f 3` != "configured" ]; then
                echo "skipping $prefix$i, already installed"
            else
                echo installing $prefix$i
                mkdir -pm 0700 $dir/$prefix$i
                chmod 700 $dir/$prefix$i
                zoneadm -z $prefix$i install > /dev/null 2>&1 &
                sleep 1 # spread things out just a tad
            fi
        fi
        i=`expr $i + 1`
        j=`expr $j + 1`
    done
    wait
done

i=1
while [ $i -le $nzones ]; do

```

```

echo setting up sysid for $prefix$i
cfg=$dir/$prefix$i/root/etc/sysidcfg
rm -f $cfg
echo "network_interface=NONE {hostname=$prefix$i}" > $cfg
echo "system_locale=C" >> $cfg
echo "terminal=xterms" >> $cfg
echo "security_policy=NONE" >> $cfg
echo "name_service=NONE" >> $cfg
echo "timezone=US/Pacific" >> $cfg
echo "root_password=Qexr7Y/wzkSbc" >> $cfg # 'lla'
i=`expr $i + 1`
done

i=1
para=`expr $nprocs \* 2`
while [ $i -le $nzones ]; do
    date
    j=1
    while [ $j -le $para ]; do
        if [ $i -le $nzones ]; then
            echo booting $prefix$i
            zoneadm -z $prefix$i boot &
        fi
        j=`expr $j + 1`
        i=`expr $i + 1`
    done
    wait
done

```

Using the zonecfg Command to Modify a Zone Configuration

You can also use the `zonecfg` command to do the following:

- Modify a resource type in a zone configuration
- Add a dedicated device to a zone

▼ How to Modify a Resource Type in a Zone Configuration

You can select a resource type and modify the specification for that resource.

Note that the contents of software packages in the `inherit-pkg-dir` directory cannot be modified or removed after the zone has been installed with `zoneadm`.

You must be the global administrator in the global zone to perform this procedure.

1. Become superuser, or assume the Primary Administrator role.

To create the role and assign the role to a user, see “Using the Solaris Management Tools With RBAC (Task Map)” in *System Administration Guide: Basic Administration*.

2. Select the zone to be modified, my-zone in this procedure.

```
global# zonecfg -z my-zone
```

3. Select the resource type to be changed, for example, a resource control.

```
zonecfg:my-zone> select rctl name=zone.cpu-shares
```

4. Remove the current value.

```
zonecfg:my-zone:rctl> remove value (priv=privileged,limit=20,action=none)
```

5. Add the new value.

```
zonecfg:my-zone:rctl> add value (priv=privileged,limit=10,action=none)
```

6. End the revised rctl specification.

```
zonecfg:my-zone:rctl> end
```

7. Commit the zone configuration for the zone.

```
zonecfg:my-zone> commit
```

8. Exit the zonecfg command.

```
zonecfg:my-zone> exit
```

Note that even if you did not explicitly type `commit` at the prompt, a `commit` is automatically attempted when you type `exit` or an EOF occurs.

▼ How to Add a Dedicated Device to a Zone

The following specification places a scanning device in a non-global zone configuration.

You must be the global administrator in the global zone to perform this procedure.

1. Become superuser, or assume the Primary Administrator role.

To create the role and assign the role to a user, see “Using the Solaris Management Tools With RBAC (Task Map)” in *System Administration Guide: Basic Administration*.

2. Add a device.

```
zonecfg:myzone> add device
```

3. Set the device match, `/dev/scsi/scanner/c3t4*` in this procedure.

```
zonecfg:myzone:device> set match=/dev/scsi/scanner/c3t4*
```

4. End the device specification.

```
zonecfg:myzone:device> end
```

5. Exit the `zonecfg` command.

```
zonecfg:myzone> exit
```

Using the `zonecfg` Command to Revert or Remove a Zone Configuration

Use the `zonecfg` command described in `zonecfg(1M)` to revert a zone's configuration or to delete a zone configuration.

▼ How to Revert a Zone Configuration

While configuring a zone with the `zonecfg` utility, use the `revert` subcommand to undo a resource setting made to the zone configuration.

You must be the global administrator in the global zone to perform this procedure.

1. Become superuser, or assume the Primary Administrator role.

To create the role and assign the role to a user, see "Using the Solaris Management Tools With RBAC (Task Map)" in *System Administration Guide: Basic Administration*.

2. While configuring a zone called `tmp-zone`, type `info` to view your configuration:

```
zonecfg:tmp-zone> info
```

The `net` resource segment of the configuration displays as follows:

```
.  
. .  
fs:  
    dir: /tmp  
    special: swap  
    type: tmpfs  
net:  
    address: 192.168.0.1  
    physical: eri0
```

```
device
    match: /dev/pts/*
.
.
.
```

3. Remove the net address:

```
zonecfg:tmp-zone> remove net address=192.168.0.1
```

4. Verify that the net entry has been removed.

```
zonecfg:tmp-zone> info
.
.
.
fs:
    dir: /tmp
    special: swap
    type: tmpfs
device
    match: /dev/pts/*
.
.
.
```

5. Type revert.

```
zonecfg:tmp-zone> revert
```

6. Answer yes to the following question:

```
Are you sure you want to revert (y/[n])? y
```

7. Verify that the net address is once again present:

```
zonecfg:tmp-zone> info
.
.
.
fs:
    dir: /tmp
    special: swap
    type: tmpfs
net:
    address: 192.168.0.1
    physical: eri0
device
    match: /dev/pts/*
.
.
.
```

▼ How to Delete a Zone Configuration

Use `zoncfg` with the `delete` subcommand to delete a zone configuration from the system.

You must be the global administrator in the global zone to perform this procedure.

- 1. Become superuser, or assume the Primary Administrator role.**

To create the role and assign the role to a user, see “Using the Solaris Management Tools With RBAC (Task Map)” in *System Administration Guide: Basic Administration*.

- 2. Delete the zone configuration for the zone `a-zone` by using one of the following two methods:**

- Use the `-F` option to force the action:

```
global# zoncfg -z a-zone delete -F
```

- Delete the zone interactively by answering yes to the system prompt:

```
global# zoncfg -z a-zone delete
Are you sure you want to delete zone a-zone (y/[n])? y
```

About Installing, Halting, and Uninstalling Non-Global Zones (Overview)

This chapter discusses zone installation on your Solaris system. It also describes the two processes that manage the virtual platform and the application environment, `zoneadm` and `zsched`. Information about halting, rebooting, and uninstalling zones is also provided.

The following topics are addressed in this chapter:

- “Zone Installation Concepts” on page 237
- “Zone Construction” on page 238
- “The `zoneadm` Daemon” on page 239
- “The `zsched` Zone Scheduler” on page 240
- “Zone Application Environment” on page 240
- “About Halting, Rebooting, and Uninstalling Zones” on page 240

To install and boot a non-global zone, or to halt or uninstall a non-global zone, see Chapter 20.

Zone Installation Concepts

The `zoneadm` command described in the `zoneadm(1M)` man page is the primary tool used to install and administer non-global zones. Operations using the `zoneadm` command must be run from the global zone. The following tasks can be performed using the `zoneadm` command:

- Verify a zone
- Install a zone
- Boot a zone
- Display information about a running zone
- Halt a zone
- Reboot a zone

- Uninstall a zone

For zone installation and verification procedures, see [Chapter 20](#) and the `zoneadm(1M)` man page. For zone configuration procedures, see [Chapter 18](#) and the `zonecfg(1M)` man page. Zone states are described in “[Non-Global Zone State Model](#)” on page 201.

If you plan to produce Solaris auditing records for zones, read “[Using Solaris Auditing in Zones](#)” on page 308 before you install non-global zones.

Zone Construction

After you have configured a non-global zone, you should verify that the zone can be installed safely on your system’s configuration. You can then install the zone. The files needed for the zone’s root file system are installed by the system under the zone’s root path. A successfully installed zone is ready for initial login and booting.

The method used to initially install packages in a Solaris installation is also the method used to populate a non-global zone.

The global zone must contain all the data necessary to populate a non-global zone. Populating a zone includes creating directories, copying files, and providing configuration information.

Only the information or data that was created in the global zone from packages is used to populate the zone from the global zone. For more information, see the `pkgparam(1)` and `pkginfo(4)` man pages.

Data from the following are not referenced or copied when a zone is installed:

- Non-installed packages
- Patches
- Data on CDs and DVDs
- Network installation images
- Any prototype or other instance of a zone

In addition, the following types of information, if present in the global zone, are not copied into a zone that is being installed:

- New or changed users in the `/etc/passwd` file
- New or changed groups in the `/etc/group` file
- Configurations for networking services such as DHCP address assignment, UUCP, or sendmail
- Configurations for network services such as naming services
- New or changed `crontab`, printer, and mail files

- System log, message, and accounting files

If Solaris auditing is used, modifications to auditing files copied from the global zone might be required. For more information, see [“Using Solaris Auditing in Zones”](#) on page 308.

The following features cannot be configured in a non-global zone:

- Solaris Live Upgrade™ boot environments
- Solaris Volume Manager metadevices
- DHCP address assignment

The resources specified in the configuration file are added when the zone transitions from installed to ready. A unique zone ID is assigned by the system. File systems are mounted, network interfaces are plumbed, and devices are configured. Transitioning into the ready state prepares the virtual platform to begin running user processes. In the ready state, the `zsched` and `zoneadmd` processes are started to manage the virtual platform.

- `zsched`, a system scheduling process similar to `sched`, is used to track kernel resources associated with the zone.
- `zoneadmd` is the zones administration daemon.

A zone in the ready state does not have any user processes executing in it. The primary difference between a ready zone and a running zone is that at least one process is executing in a running zone. See the `init(1M)` man page for more information.

The zoneadmd Daemon

The zones administration daemon, `zoneadmd`, is the primary process for managing the zone’s virtual platform. The daemon is also responsible for managing zone booting and shutting down. There is one `zoneadmd` process running for each active (ready, running, or shutting down) zone on the system.

The `zoneadmd` daemon sets up the zone as specified in the zone configuration. This process includes the following actions:

- Allocating the zone ID and starting the `zsched` system process.
- Setting zone-wide resource controls.
- Preparing the zone’s devices as specified in the zone configuration. For more information, see the `devfsadmd(1M)` man page.
- Plumbing virtual network interfaces.
- Mounting loopback and conventional file systems.

- Instantiating and initializing the zone console device.

Unless the `zoneadm` daemon is already running, it is automatically started by `zoneadm`. Thus, if the daemon is not running for any reason, any invocation of `zoneadm` to administer the zone will restart `zoneadm`.

The man page for the `zoneadm` daemon is `zoneadm(1M)`.

The `zsched` Zone Scheduler

An active zone is a zone that is in the ready state, the running state, or the shutting down state. Every active zone has an associated kernel process, `zsched`. Kernel threads doing work on behalf of the zone are owned by `zsched`. The `zsched` process enables the zones subsystem to keep track of per-zone kernel threads.

Zone Application Environment

Booting a zone is similar to booting a regular Solaris system. The `zoneadm` command is used to create the zone application environment.

Before a non-global zone is booted for the first time, the internal configuration of the zone must be created. The internal configuration specifies a naming service to use, the default locale and time zone, the zone's root password, and other aspects of the application environment. The application environment is established by responding to a series of prompts that appear on the zone console, as explained in [“Internal Zone Configuration” on page 254](#). Note that the default locale and time zone for a zone can be configured independently of the global settings.

About Halting, Rebooting, and Uninstalling Zones

This section provides an overview of the procedures for halting, rebooting, and uninstalling zones. Troubleshooting tips for zones that fail to halt when requested are also provided.

Halting a Zone

The `zoneadm halt` command is used to remove both the application environment and the virtual platform for a zone. The zone is then brought back to the installed state. All processes are killed, devices are unconfigured, network interfaces are unplumbed, file systems are unmounted, and the kernel data structures are destroyed.

The `halt` command does *not* run any shutdown scripts within the zone. To shut down a zone, see [“How to Use `zlogin` to Shut Down a Zone”](#) on page 264.

In the event that the system state associated with the zone cannot be destroyed, the `halt` operation will fail halfway. This leaves the zone in an intermediate state, somewhere between running and installed. In this state there are no active user processes or kernel threads, and none can be created. When the `halt` operation fails, you must manually intervene to complete the process.

The most common cause of a failure is the inability of the system to unmount all file systems. Unlike a traditional Solaris system shutdown, which destroys the system state, zones must ensure that no mounts performed while booting the zone or during zone operation remain once the zone has been halted. Even though `zoneadm` makes sure that there are no processes executing in the zone, the unmount operation can fail if processes in the global zone have open files in the zone. Use the tools described in the `proc(1)` (see `pfiles`) and `fuser(1M)` man pages to find these processes and take appropriate action. After these processes have been dealt with, reinvoking `zoneadm halt` will completely halt of the zone.

Rebooting a Zone

The `zoneadm reboot` command is used to reboot a zone. The zone is halted and then booted again. The zone ID will change when the zone is rebooted.

Zone autoboot

If you set the `autoboot` resource property in a zone's configuration to `true`, that zone is automatically booted when the global zone is booted. The default setting is `false`.

Uninstalling a Zone

The `zoneadm uninstall` command is used to uninstall all of the files under the zone's root file system. Before proceeding, the command prompts you to confirm the action, unless the `-F` (force) option is also used. Use the `uninstall` command with caution, because the action is irreversible.

Installing, Booting, Halting, and Uninstalling Non-Global Zones (Tasks)

This chapter describes how to install and boot a non-global zone. Other tasks associated with installation, such as halting, rebooting, and uninstalling zones, are also addressed. The procedure to completely delete a zone from a system is also provided.

For general information about zone installation and related operations, see [Chapter 19](#).

Zone Installation (Task Map)

Task	Description	For Instructions
(Optional) Verify a configured zone prior to installing the zone.	Ensure that a zone meets the requirements for installation. If you skip this procedure, the verification is performed automatically when you install the zone.	“How to Verify a Configured Zone Before It Is Installed (Optional)” on page 244
Install a configured zone.	Install a zone that is in the configured state.	“How to Install a Configured Zone” on page 245
(Optional) Transition an installed zone to the ready state.	You can skip this procedure if you want to boot the zone and use it immediately.	“How to Transition the Installed Zone to the Ready State (Optional)” on page 246

Task	Description	For Instructions
Boot a zone.	Booting a zone places the zone in the running state. A zone can be booted from the ready state or from the installed state. Note that you must perform the internal zone configuration before you boot the zone for the first time. This is described in “Internal Zone Configuration” on page 254.	“How to Boot a Zone” on page 246, “Performing the Initial Internal Zone Configuration” on page 258

Installing and Booting Zones

Use the `zoneadm` command described in the `zoneadm(1M)` man page to perform installation tasks for a non-global zone. You must be the global administrator to perform the zone installation. The examples in this chapter use the zone name and zone path established in [“Configuring, Verifying, and Committing a Zone” on page 226.](#)

▼ How to Verify a Configured Zone Before It Is Installed (Optional)

You can verify a zone prior to installing it. If you skip this procedure, the verification is performed automatically when you install the zone.

You must be the global administrator in the global zone to perform this procedure.

1. Become superuser, or assume the Primary Administrator role.

To create the role and assign the role to a user, see [“Using the Solaris Management Tools With RBAC \(Task Map\)” in *System Administration Guide: Basic Administration.*](#)

2. Verify a configured zone named `my-zone` by using the `-z` option with the name of the zone and the `verify` subcommand.

```
global# zoneadm -z my-zone verify
```

If an error message is displayed and the zone fails to verify, make the corrections specified in the message and try the command again.

For example, `zoneadm verify` might issue a warning if the zone path does not exist because of a typing error:

Warning: /export/home1/my-zone does not exist, so it cannot be verified.
When 'zoneadm install' is run, 'install' will try to create
/export/home1/my-zone, and 'verify' will be tried again,
but the 'verify' may fail if:
the parent directory of /export/home1/my-zone is group- or other-writable
or
/export/home1/my-zone overlaps with any other installed zones.

If no error messages are displayed, you can install the zone.

▼ How to Install a Configured Zone

You must be the global administrator in the global zone to perform this procedure.

1. Become superuser, or assume the Primary Administrator role.

To create the role and assign the role to a user, see “Using the Solaris Management Tools With RBAC (Task Map)” in *System Administration Guide: Basic Administration*.

2. Install the configured zone my-zone by using the zoneadm command with the -z install option.

```
global# zoneadm -z my-zone install
```

You will see various messages as the files and directories needed for the zone's root file system are installed under the zone's root path.

If an error message is displayed and the zone fails to install, type the following to get the zone state:

```
global# zoneadm -z my-zone list -v
```

- If the state is listed as configured, make the corrections specified in the message and try the zoneadm install command again.
- If the state is listed as incomplete, first execute this command:

```
global# zoneadm -z my-zone uninstall
```

Then make the corrections specified in the message, and try the zoneadm install command again.

3. When the installation completes, use the list subcommand with the -i and -v options to list the installed zones and verify the status.

```
global# zoneadm list -iv
```

You will see a display that is similar to the following:

ID	NAME	STATE	PATH
0	global	running	/
-	my-zone	installed	/export/home/my-zone

Note – If a zone installation is interrupted or fails, the zone is left in the incomplete state. Use `uninstall -F` to reset the zone to the configured state.

▼ How to Transition the Installed Zone to the Ready State (Optional)

Transitioning into the ready state prepares the virtual platform to begin running user processes. Zones in the ready state do not have any user processes executing in them.

You can skip this procedure if you want to boot the zone and use it immediately. The transition through the ready state is performed automatically when you boot the zone.

You must be the global administrator in the global zone to perform this procedure.

1. Become superuser, or assume the Primary Administrator role.

To create the role and assign the role to a user, see “Using the Solaris Management Tools With RBAC (Task Map)” in *System Administration Guide: Basic Administration*.

2. Use the `zoneadm` command with the `-z` option, the name of the zone, which is `my-zone`, and the `ready` subcommand to transition the zone to the ready state.

```
global# zoneadm -z my-zone ready
```

3. At the prompt, use the `zoneadm list` command with the `-v` option to verify the status.

```
global# zoneadm list -v
```

You will see a display that is similar to the following:

ID	NAME	STATE	PATH
0	global	running	/
1	my-zone	ready	/export/home/my-zone

Note that the unique zone ID 1 has been assigned by the system.

▼ How to Boot a Zone

Booting a zone places the zone in the running state. A zone can be booted from the ready state or from the installed state. A zone in the installed state that is booted transparently transitions through the ready state to the running state.

Note that you should log in and perform the internal zone configuration before you boot the zone for the first time. This is described in [“Internal Zone Configuration” on page 254](#).

If you plan to use an `/etc/sysidcfg` file to perform initial zone configuration, as described in “How to Use an `/etc/sysidcfg` File to Perform the Initial Zone Configuration” on page 260, create the `sysidcfg` file and place it in the zone’s `/etc` directory before you boot the zone.

You must be the global administrator in the global zone to perform this procedure.

1. Become superuser, or assume the Primary Administrator role.

To create the role and assign the role to a user, see “Using the Solaris Management Tools With RBAC (Task Map)” in *System Administration Guide: Basic Administration*.

2. Use the `zoneadm` command with the `-z` option, the name of the zone, which is `my-zone`, and the `boot` subcommand to boot the zone.

```
global# zoneadm -z my-zone boot
```

3. When the installation completes, use the `list` subcommand with the `-v` option to verify the status.

```
global# zoneadm list -v
```

You will see a display that is similar to the following:

ID	NAME	STATE	PATH
0	global	running	/
1	my-zone	running	/export/home/my-zone

If you see the following message when you boot the zone:

```
# zoneadm -z my-zone boot
zoneadm: zone 'my-zone': WARNING: hme0:1: no matching subnet
found in netmasks(4) for 192.168.0.1; using default of
255.255.255.0.
```

The message is only a warning, and the command has succeeded. The message indicates that the system was unable to find the netmask to be used for the IP address specified in the zone’s configuration.

To stop the warning from displaying on subsequent reboots, ensure that the correct `netmasks` databases are listed in the `/etc/nsswitch.conf` file in the global zone and that at least one of these databases contains the subnet and netmasks to be used for the zone `my-zone`.

For example, if the `/etc/inet/netmasks` file and the local NIS database are used for resolving netmasks in the global zone, the appropriate entry in `/etc/nsswitch.conf` is as follows:

```
netmasks: files nis
```

The subnet and corresponding netmask information for the zone `my-zone` can then be added to `/etc/inet/netmasks` for subsequent use.

For more information about the `netmasks` command, see the `netmasks(4)` man page.

Where to Go From Here

To log in to the zone and perform the initial internal configuration, see [Chapter 21](#) and [Chapter 22](#).

Halting, Rebooting, Uninstalling, and Deleting Non-Global Zones (Task Map)

Task	Description	For Instructions
Halt a zone.	The halt procedure is used to remove both the application environment and the virtual platform for a zone. The procedure returns a zone in the ready state to the installed state. To cleanly shut down a zone, see “How to Use zlogin to Shut Down a Zone” on page 264.	“How to Halt a Zone” on page 249
Reboot a zone.	The reboot procedure halts the zone and then boots it again.	“How to Reboot a Zone” on page 250
Uninstall a zone.	Removes all of the files in the zone’s root file system. <i>Use this procedure with caution.</i> The action is irreversible.	“How to Uninstall a Zone” on page 250
Delete a non-global zone from the system.	This procedure completely removes a zone from a system.	“Deleting a Non-Global Zone From the System” on page 251

Halting, Rebooting, and Uninstalling Zones

▼ How to Halt a Zone

The halt procedure is used to remove both the application environment and the virtual platform for a zone. To cleanly shut down a zone, see [“How to Use `zlogin` to Shut Down a Zone”](#) on page 264.

You must be the global administrator in the global zone to perform this procedure.

1. **Become superuser, or assume the Primary Administrator role.**

To create the role and assign the role to a user, see [“Using the Solaris Management Tools With RBAC \(Task Map\)”](#) in *System Administration Guide: Basic Administration*.

2. **List the zones running on the system.**

```
global# zoneadm list -v
```

You will see a display that is similar to the following:

ID	NAME	STATE	PATH
0	global	running	/
1	my-zone	running	/export/home/my-zone

3. **Use the `zoneadm` command with the `-z` option, the name of the zone, for example, `my-zone`, and the `halt` subcommand to halt the given zone.**

```
global# zoneadm -z my-zone halt
```

4. **List the zones on the system again, to verify that `my-zone` has been halted.**

```
global# zoneadm list -iv
```

You will see a display that is similar to the following:

ID	NAME	STATE	PATH
0	global	running	/
-	my-zone	installed	/export/home/my-zone

5. **Boot the zone if you want to restart it.**

```
global# zoneadm -z my-zone boot
```

If the zone does not halt properly, see [“Halting a Zone”](#) on page 241 for troubleshooting tips.

▼ How to Reboot a Zone

You must be the global administrator in the global zone to perform this procedure.

1. **Become superuser, or assume the Primary Administrator role.**

To create the role and assign the role to a user, see “Using the Solaris Management Tools With RBAC (Task Map)” in *System Administration Guide: Basic Administration*.

2. **List the zones running on the system.**

```
global# zoneadm list -v
```

You will see a display that is similar to the following:

ID	NAME	STATE	PATH
0	global	running	/
1	my-zone	running	/export/home/my-zone

3. **Use the zoneadm command with the -z reboot option to reboot the zone my-zone.**

```
global# zoneadm -z my-zone reboot
```

4. **List the zones on the system again to verify that my-zone has been rebooted.**

```
global# zoneadm list -v
```

You will see a display that is similar to the following:

ID	NAME	STATE	PATH
0	global	running	/
2	my-zone	running	/export/home/my-zone

Note that the zone ID for my-zone has changed. The zone ID generally changes after a reboot.

▼ How to Uninstall a Zone

Use this procedure with caution. The action of removing all of the files in the zone’s root file system is irreversible.

You must be the global administrator in the global zone to perform this procedure.

1. **Become superuser, or assume the Primary Administrator role.**

To create the role and assign the role to a user, see “Using the Solaris Management Tools With RBAC (Task Map)” in *System Administration Guide: Basic Administration*.

2. **List the zones on the system.**

```
global# zoneadm list -v
```

You will see a display that is similar to the following:

ID	NAME	STATE	PATH
0	global	running	/

```
- my-zone installed /export/home/my-zone
```

3. Use the `zoneadm` command with the `-z uninstall` option to remove the zone `my-zone`.

You can also use the `-F` option to force the action. If this option is not specified, the system will prompt for confirmation.

```
global# zoneadm -z my-zone uninstall -F
```

4. List the zones on the system again, to verify that `my-zone` is no longer listed.

```
global# zoneadm list -v
```

You will see a display that is similar to the following:

ID	NAME	STATE	PATH
0	global	running	/

Note – If a zone uninstall is interrupted, the zone is left in the incomplete state. Use the `zoneadm uninstall` command to reset the zone to the configured state. Use the `uninstall` command with caution because the action is irreversible.

Deleting a Non-Global Zone From the System

The procedure described in this section completely deletes a zone from a system.

▼ How to Remove a Non-Global Zone

1. Shut down the zone `my-zone`.

```
global# zlogin my-zone shutdown
```

2. Remove the root file system for `my-zone`.

```
global# zoneadm -z my-zone uninstall -F
```

3. Delete the configuration for `my-zone`.

```
global# zonecfg -z my-zone delete -F
```

4. List the zones on the system, to verify that `my-zone` is no longer listed.

```
global# zoneadm list -iv
```

You will see a display that is similar to the following:

ID	NAME	STATE	PATH
0	global	running	/

Non-Global Zone Login (Overview)

This chapter discusses logging in to zones from the global zone.

The following topics are covered in this chapter:

- “zlogin Command” on page 253
- “Non-Global Zone Login Methods” on page 255
- “Interactive and Non-Interactive Modes” on page 256
- “Failsafe Mode” on page 255
- “Remote Login” on page 256

For procedures and usage information, see [Chapter 22](#).

zlogin Command

After you install a zone, you must log in to the zone to complete its application environment. You might log in to the zone to perform administrative tasks as well. Unless the `-C` option is used to connect to the zone console, logging in to a zone using `zlogin` starts a new task. A task cannot span two zones.

The `zlogin` command is used to log in from the global zone to any zone that is in the running state or the ready state.

Note – Only the `zlogin` command with the `-C` option can be used to log in to a zone that is not in the running state.

As described in [“How to Use Non-Interactive Mode to Access a Zone”](#) on page 263, you can use the `zlogin` command in non-interactive mode by supplying a command to run inside a zone. However, the command or any files the command acts upon cannot reside on NFS. The command will fail if any of its open files or any portion of its address space resides on NFS. The address space includes the command executable itself and the command’s linked libraries.

The `zlogin` command can only be used by the global administrator operating in the global zone. See the `zlogin(1)` man page for more information.

Internal Zone Configuration

When a zone is booted for the first time after installation, the zone is in an unconfigured state. The zone does not have an internal configuration for naming services, its locale and time zone have not been set, and various other configuration tasks have not been performed. Therefore, the `sysidtool` programs are run the first time a zone is booted. For more information, see the `sysidtool(1M)` man page.

Two methods are available for performing the required configuration:

- Zone console login, which initiates a series of questions from the system. Be prepared to respond to the following:
 - Language
 - Type of terminal being used
 - Host name
 - Security policy (Kerberos or standard UNIX)
 - Naming service type (None is a valid response)
 - Naming service domain
 - Name server
 - Default time zone
 - Root password

The procedure is described in [“Performing the Initial Internal Zone Configuration”](#) on page 258.

- An `/etc/sysidcfg` file, which you can create and place inside the zone before you boot the zone for the first time. See the `sysidcfg(4)` man page for more information.

Non-Global Zone Login Methods

This section describes the methods you can use to log in to a zone.

Zone Console Login

Each zone maintains a virtual console, `/dev/console`. Performing actions on the console is referred to as console mode. Connections to the console persist across zone reboots.

The zone console is accessed by using the `zlogin` command with the `-C` option and the *zonename*. The zone does not have to be in the running state.

Processes inside the zone can open and write messages to the console. If the `zlogin -C` process exits, another process can then access the console.

User Login Methods

To log in to the zone with a user name, use the `zlogin` command with the `-l` option, the user name, and the *zonename*. For example, the administrator of the global zone can log in as a normal user in the non-global zone by specifying the `-l` option to `zlogin`:

```
global# zlogin -l user zonename
```

To log in as user `root`, use the `zlogin` command without options.

Failsafe Mode

If a login problem occurs and you cannot use the `zlogin` command or the `zlogin` command with the `-C` option to access the zone, an alternative is provided. You can enter the zone by using the `zlogin` command with the `-S` (safe) option. Only use this mode to recover a damaged zone when other forms of login are not succeeding. In this minimal environment, it might be possible to diagnose why the zone login is failing.

Remote Login

The ability to remotely log in to a zone is dependent on the selection of network services that you establish. By default, logins through `rlogin`, `ssh`, and `telnet` function normally. For more information about these commands, see `rlogin(1)`, `ssh(1)`, and `telnet(1)`.

Interactive and Non-Interactive Modes

Two other methods for accessing the zone and for executing commands inside the zone are also provided by the `zlogin` command. These methods are interactive mode and non-interactive mode.

Interactive Mode

In interactive mode, a new pseudo-terminal is allocated for use inside the zone. Unlike console mode, in which exclusive access to the console device is granted, an arbitrary number of `zlogin` sessions can be open at any time in interactive mode. Interactive mode is activated when you do not include a command to be issued. Programs that require a terminal device, such as an editor, operate correctly in this mode.

Non-Interactive Mode

Non-interactive mode is used to run shell-scripts which administer the zone. Non-interactive mode does not allocate a new pseudo-terminal. Non-interactive mode is enabled when you supply a command to be run inside the zone.

Logging In to Non-Global Zones (Tasks)

This chapter provides procedures for completing the configuration of an installed zone, logging into a zone from the global zone, and shutting down a zone. This chapter also shows how to use the `zonename` command to print the name of the current zone

For an introduction to the zone login process, see [Chapter 21](#).

Initial Zone Boot and Zone Login Procedures (Task Map)

Task	Description	For Instructions
Perform the internal configuration.	Log in to the zone console or use an <code>/etc/sysidcfg</code> file to perform the initial zone configuration.	“Performing the Initial Internal Zone Configuration” on page 258

Task	Description	For Instructions
Log in to the zone.	You can log into a zone through the console, by using interactive mode to allocate a pseudo-terminal, or by supplying a command to be run in the zone. Supplying a command to be run does not allocate a pseudo-terminal. You can also log in by using failsafe mode when a connection to the zone is denied.	“Logging In to a Zone” on page 261
Shut down a zone.	Shut down a zone by using the <code>shutdown</code> utility or a script.	“How to Use <code>zlogin</code> to Shut Down a Zone” on page 264
Print the zone name.	Print the zone name of the current zone.	“Printing the Name of the Current Zone” on page 264

Performing the Initial Internal Zone Configuration

Before you boot a zone after installation, you must configure the zone using one of the following methods:

- Log into the zone and configure it as described in [“Internal Zone Configuration” on page 254](#).
- Configure the zone using an `/etc/sysidcfg` file as described in [“How to Use an `/etc/sysidcfg` File to Perform the Initial Zone Configuration” on page 260](#).

▼ How to Log In to the Zone Console to Perform the Internal Zone Configuration

You must be the global administrator in the global zone to perform this procedure.

1. Become superuser, or assume the Primary Administrator role.

To create the role and assign the role to a user, see [“Using the Solaris Management Tools With RBAC \(Task Map\)”](#) in *System Administration Guide: Basic Administration*.

2. Use the `zlogin` command with the `-C` option and the name of the zone, `my-zone` in this procedure.

```
global# zlogin -C my-zone
```

3. The first time you log in to the console, you are prompted to answer a series of questions. Your screen will look similar to this:

```
SunOS Release 5.10 Version Generic 64-bit
Copyright 1983-2004 Sun Microsystems, Inc. All rights reserved.
Use is subject to license terms.
```

```
Select a Language
```

- 0. English
- 1. French

```
Please make a choice (0 - 1), or press h or ? for help:
```

```
Select a Locale
```

- 0. English (C - 7-bit ASCII)
- 1. Belgium-Flemish (ISO8859-1)
- 2. Belgium-Flemish (ISO8859-15 - Euro)
- 3. Great Britain (ISO8859-1)
- 4. Great Britain (ISO8859-15 - Euro)
- 5. Ireland (ISO8859-1)
- 6. Ireland (ISO8859-15 - Euro)
- 7. Netherlands (ISO8859-1)
- 8. Netherlands (ISO8859-15 - Euro)
- 9. Go Back to Previous Screen

```
Please make a choice (0 - 9), or press h or ? for help:
```

```
What type of terminal are you using?
```

- 1) ANSI Standard CRT
- 2) DEC VT52
- 3) DEC VT100
- 4) Heathkit 19
- 5) Lear Siegler ADM31
- 6) PC Console
- 7) Sun Command Tool
- 8) Sun Workstation
- 9) Televideo 910
- 10) Televideo 925
- 11) Wyse Model 50
- 12) X Terminal Emulator (xterms)
- 13) CDE Terminal Emulator (dtterm)
- 14) Other

```
Type the number of your choice and press Return:
```

```
.
.
.
```

For the complete list of questions you must answer, see [“Internal Zone Configuration”](#) on page 254.

If you booted this zone before you logged in, you might have missed the initial prompt for configuration information. If you see the following system message at zone login instead of a prompt:

```
[connected to zone zonename console]
```

Press Return to display the prompt again.

If you enter an incorrect response and try to restart the configuration, you might experience difficulty when you attempt the process again. This occurs because the `sysidtools` can store your previous responses.

If this happens, use the following workaround from the global zone to restart the configuration process.

```
global# zlogin -S zonename /usr/sbin/sys-unconfig
```

For more information on the `sys-unconfig` command, see the `sys-unconfig(1M)` man page.

▼ How to Use an `/etc/sysidcfg` File to Perform the Initial Zone Configuration

You must be the global administrator in the global zone to perform this procedure.

1. Become superuser, or assume the Primary Administrator role.

To create the role and assign the role to a user, see “Using the Solaris Management Tools With RBAC (Task Map)” in *System Administration Guide: Basic Administration*.

2. From the global zone, change directories to the non-global zone’s `/etc` directory:

```
global# cd /export/home/my-zone/root/etc
```

3. Create the `sysidcfg` file and place it in this directory.

The file will look similar to the following:

```
system_locale=C
terminal=dtterm
network_interface=primary {
    hostname=my-zone
}
security_policy=NONE
name_service=NIS {
    domain_name=special.example.com
    name_server=bird(192.168.112.3)
}
timezone=US/Central
root_password=m4qt0WN
```

4. By default, a separate module will request the NFSv4 domain parameter used by the `nfsmapid` command. To complete a hands-off initial zone configuration, edit

the file `default/nfs`, uncomment the `NFSMAPID_DOMAIN` parameter, and set the domain to the desired NFSv4 domain:

```
global# vi default/nfs
.
.
.
NFSMAPID_DOMAIN=domain
```

For more information on the NFSv4 domain parameter, see the `nfsmapid(1M)` man page.

5. Create the file `.NFS4inst_state.domain` in this directory to indicate that the NFSv4 domain has been set:

```
global# touch .NFS4inst_state.domain
```

6. Boot the zone.

Logging In to a Zone

Use the `zlogin` command to log in from the global zone to any zone that is running or in the ready state. See the `zlogin(1)` man page for more information.

You can log in to a zone in various ways, as described in the following procedures. You can also log in remotely, as described in [“Remote Login” on page 256](#).

▼ How to Log In to the Zone Console

You must be the global administrator in the global zone to perform this procedure.

1. **Become superuser, or assume the Primary Administrator role.**

To create the role and assign the role to a user, see *“Using the Solaris Management Tools With RBAC (Task Map)”* in *System Administration Guide: Basic Administration*.

2. **Use the `zlogin` command with the `-C` option and the name of the zone, for example, `my-zone`.**

```
global# zlogin -C my-zone
```

Note – If you start the `zlogin` session immediately after issuing the `zoneadm boot` command, boot messages from the zone will display:

```
SunOS Release 5.10 Version Generic 64-bit
Copyright 1983-2004 Sun Microsystems, Inc. All rights reserved.
Use is subject to license terms.
starting rpc services: rpcbind done.
syslog service starting.
The system is ready.
```

3. **When the zone console displays, log in as root, press Return, and type the root password when prompted.**

```
my-zone console login: root
Password:
```

▼ How to Use Interactive Mode to Access a Zone

In interactive mode, a new pseudo-terminal is allocated for use inside the zone.

You must be the global administrator in the global zone to perform this procedure.

1. **Become superuser, or assume the Primary Administrator role.**

To create the role and assign the role to a user, see “Using the Solaris Management Tools With RBAC (Task Map)” in *System Administration Guide: Basic Administration*.

2. **At the global zone prompt, type `tty`.**

```
global# tty
The pseudo-terminal device name will display:
/dev/pts/3
```

3. **From the global zone, log in to the zone, for example, `my-zone`.**

```
global# zlogin my-zone
Information similar to the following will display:
[Connected to zone 'my-zone' pts/2]
Last login: Wed Jul  3 16:25:00 on console
Sun Microsystems Inc. SunOS 5.10 Generic July 2004
```

4. **At the `my-zone` prompt, type `tty`.**

```
my-zone# tty
The name of the pseudo-terminal device displays:
/dev/pts/2
my-zone#
```

5. Type `exit` to close the connection.

You will see a messages similar to the following:

```
[Connection to zone 'my-zone' pts/2 closed]
```

▼ How to Use Non-Interactive Mode to Access a Zone

Non-interactive mode is enabled when the user supplies a command to be run inside the zone. Non-interactive mode does not allocate a new pseudo-terminal.

Note that the command or any files that the command acts upon cannot reside on NFS.

You must be the global administrator in the global zone to perform this procedure.

1. **Become superuser, or assume the Primary Administrator role.**

To create the role and assign the role to a user, see “Using the Solaris Management Tools With RBAC (Task Map)” in *System Administration Guide: Basic Administration*.

2. **From the global zone, log in to the `my-zone` zone and supply a command name.**

The command `zonename` is used here.

```
global# zlogin my-zone zonename
```

You will see the following output:

```
my-zone
```

▼ How to Use Failsafe Mode to Enter a Zone

When connection to the zone is denied, the `zlogin` command can be used with the `-S` option to enter a minimal environment in the zone.

You must be the global administrator in the global zone to perform this procedure.

1. **Become superuser, or assume the Primary Administrator role.**

To create the role and assign the role to a user, see “Using the Solaris Management Tools With RBAC (Task Map)” in *System Administration Guide: Basic Administration*.

2. **From the global zone, use the `zlogin` command with the `-S` option to access the zone, for example, `my-zone`.**

```
global# zlogin -S my-zone
```

▼ How to Use `zlogin` to Shut Down a Zone

Use this procedure to cleanly shut down a zone. To halt a zone without running shutdown scripts, see “How to Halt a Zone” on page 249.

You must be the global administrator in the global zone to perform this procedure.

1. **Become superuser, or assume the Primary Administrator role.**

To create the role and assign the role to a user, see “Using the Solaris Management Tools With RBAC (Task Map)” in *System Administration Guide: Basic Administration*.

2. **Log in to the zone to be shut down, for example, `my-zone`. Specify `shutdown` as the name of the utility.**

```
global# zlogin my-zone shutdown
```

Your site might have its own shutdown script, tailored for your specific environment.

Note – Running `init 0` in the global zone to cleanly shut down a Solaris system also runs `init 0` in each of the non-global zones on the system. Note that `init 0` does not warn local and remote users to log off before the system is taken down.

Printing the Name of the Current Zone

The `zonename` command described in the `zonename(1)` man page prints the name of the current zone. The following example shows the output when `zonename` is used in the global zone.

```
# zonename
global
```


About Adding and Removing Packages and Patches in a Solaris Zones Environment (Overview)

This chapter discusses maintaining the Solaris Operating System in a zones environment. Information about adding packages and patches to the operating system in the global zone and in all installed non-global zones is provided. Information about removing packages and patches is also included. The material in this chapter supplements the existing Solaris installation and patch documentation. See the *Solaris 10 Release and Installation Collection* and *System Administration Guide: Basic Administration* for more information.

This chapter covers the following topics:

- “Packaging and Patch Tools Overview” on page 265
- “About Packages and Zones” on page 267
- “Keeping Zones in Sync” on page 267
- “About Adding Packages in Zones” on page 269
- “About Removing Packages in Zones” on page 272
- “Package Parameter Information” on page 274
- “Package Information Query” on page 278
- “About Adding Patches in Zones” on page 279
- “Applying Patches in a Zones Environment” on page 279
- “Removing Patches in a Zones Environment” on page 280
- “PatchPro Support” on page 281
- “Product Database” on page 281

Packaging and Patch Tools Overview

The Solaris packaging tools are used in administering the zones environment. The global administrator can upgrade the system to a new version of Solaris, which updates both the global and the non-global zones.

Note – Solaris Live Upgrade software can be used in the global zone to upgrade a system that includes non-global zones. No other upgrade method is supported in a zones environment.

The zone administrator can use the packaging tools to administer any software installed in a non-global zone, within the limits described in this document.

The following general principles apply in a zones environment:

- The global administrator can administer the software on every zone on the system.
- The root file system for a non-global zone can be administered from the global zone by using the Solaris packaging and patch tools. The Solaris packaging and patch tools are supported within the non-global zone for administering bundled, unbundled, or third-party products.
- The packaging and patch tools work in a zones-enabled environment. The tools allow a package or patch installed in the global zone to also be installed in a non-global zone.
- The `SUNW_PKG_ALLZONES` package parameter defines the *zone scope* of a package. The scope determines the type of zone in which an individual package can be installed. For more information about this parameter, see [“SUNW_PKG_ALLZONES Package Parameter” on page 274](#).
- The `SUNW_PKG_HOLLOW` package parameter defines the *visibility* of a package if that package is required to be installed on all zones and be identical in all zones. For information about this parameter, see [“SUNW_PKG_HOLLOW Package Parameter” on page 276](#).
- The packaging information visible from within a non-global zone is consistent with the files that have been installed in that zone using the Solaris packaging and patch tools. The visibility includes packages that have been imported from the global zone using read-only loopback mounts. See [“Configuring, Verifying, and Committing a Zone” on page 226](#) for more information about this process.
- A change, such as a Solaris upgrade or a patch added in the global zone, can be pushed out to all of the zones. This feature maintains consistency between the global zone and each non-global zone.
- The package commands can add, remove, and interrogate packages. The patch commands can add and remove patches.

Note – While operations are performed, a zone is temporarily locked. The system will also confirm a requested operation with the administrator before proceeding.

About Packages and Zones

Only a subset of the Solaris packages installed on the global zone are completely replicated when a non-global zone is installed. For example, many packages that contain the Solaris kernel are not needed in a non-global zone. All non-global zones implicitly share the same Solaris kernel from the global zone. However, even if a package's data is not required or is not of use in a non-global zone, the knowledge that a package is installed in the global zone might be required in a non-global zone. The information allows package dependencies from the non-global zones to be properly resolved with the global zone.

Packages have parameters that control how their content is distributed and made visible on a system with non-global zones installed. The `SUNW_PKG_ALLZONES` and `SUNW_PKG_HOLLOW` package parameters define the characteristics of packages in a zones environment. If desired, system administrators can check these package parameter settings to verify the patch's applicability when applying or removing a package in a zone environment. The `pkgparam` command can be used to view the values for these parameters. For more information on parameters, see [“Package Parameter Information” on page 274](#). See [“Checking Package Parameter Settings in a Zones Environment” on page 291](#) for usage instructions.

For information about package characteristics and parameters, see the `pkginfo(4)` man page. For information about displaying package parameter values, see the `pkgparam(1)` man page.

Note – Any package that must be interactive, which means that it has a request script, is added to the current zone only. If an interactive package is added to the global zone, the package is treated as though it is being added by using the `pkgadd` command with the `-G` option. For more information about this option, see [“About Adding Packages in Zones” on page 269](#).

Keeping Zones in Sync

It is best to keep the software installed in the non-global zones in sync with the software installed in the global zone to the maximum extent possible. This practice minimizes the difficulty in administering a system with multiple installed zones.

To achieve this goal, the package tools enforce the following rules when adding or removing packages in the global zone.

Package Operations Possible in the Global Zone

If the package is not currently installed in the global zone and not currently installed in any non-global zone, the package can be installed:

- Only in the global zone, if `SUNW_PKG_ALLZONES=false`
- In the global zone and all non-global zones

If the package is currently installed in the global zone only:

- The package can be installed in all non-global zones.
- The package can be removed from the global zone.

If a package is currently installed in the global zone and currently installed in only a subset of the non-global zones:

- `SUNW_PKG_ALLZONES` must be set to `false`.
- The package can be installed in all non-global zones. Existing instances in any non-global zone are updated to the revision being installed.
- The package can be removed from the global zone.
- The package can be removed from the global zone and from all non-global zones.

If a package is currently installed in the global zone and currently installed in all non-global zones, the package can be removed from the global zone and from all non-global zones.

These rules ensure the following:

- Packages installed in the global zone are either installed in the global zone only, or installed in the global zone and all non-global zones.
- Packages installed in the global zone and also installed in any non-global zone are the same across all zones.

Package Operations Possible in a Non-Global Zone

The package operations possible in any non-global zone are:

- If a package is not currently installed in the non-global zone, the package can be installed only if `SUNW_PKG_ALLZONES=false`.
- If a package is currently installed in the non-global zone:
 - The package can be installed over the existing instance of the package only if `SUNW_PKG_ALLZONES=false`.
 - The package can be removed from the non-global zone only if `SUNW_PKG_ALLZONES=false`.

About Adding Packages in Zones

The `pkgadd` system utility described in the `pkgadd(1M)` man page is used to add packages in a zones environment.

Using `pkgadd` in the Global Zone

The following `pkgadd` utility options are used from the global zone:

- G Add package(s) to the global zone only.
- Z Add package(s) already installed in the global zone to the non-global zones only.

When you run the `pkgadd` utility in the global zone, the following actions apply:

- The `pkgadd` utility is able to add a package:
 - To the global zone only, unless the package is `SUNW_PKG_ALLZONES=true`
 - To the global zone and to all non-global zones
 - To all non-global zones only, if the package is already installed in the global zone
- The `pkgadd` utility cannot add a package:
 - To any subset of the non-global zones
 - To all non-global zones, unless the package is already installed in the global zone
- If the `pkgadd` utility is run without the `-G` option and the `-Z` option, the specified package is added to all zones by default. The package is not marked as installed in the global zone only.
- If the `-G` option is used, the `pkgadd` utility adds the specified package to the global zone only. The package is marked as installed in the global zone only. The package is not installed when any non-global zone is installed.
- If the `-Z` option is used, the `pkgadd` utility adds the specified package to all non-global zones only. The package must already be installed in the global zone. The package is not marked as installed in the global zone only.

Adding a Package to the Global Zone and to all Non-Global Zones

To add a package to the global zone and to all non-global zones, execute the `pkgadd` utility in the global zone. As the global administrator, run `pkgadd` without the `-G` option or the `-Z` option.

A package can be added to the global zone and to all non-global zones without regard to the area affected by the package.

The following steps are performed by the `pkgadd` utility:

- Package dependencies are checked on the global zone and on all non-global zones. If required packages are not installed in any zone, then the dependency check fails. The system notifies the global administrator, who is prompted whether to continue.
- The package is added to the global zone.
- The package database on the global zone is updated.
- The package is added to each non-global zone and the database in the global zone is updated.
- The package database on each non-global zone is updated.

Adding a Package to the Global Zone Only

To add a package to the global zone only, as the global administrator in the global zone, execute the `pkgadd` utility with the `-G` option only. Do not use the `-Z` option.

A package can be added to the global zone if the following conditions are true:

- The package contents do not affect any area of the global zone that is shared with any non-global zone.
- The package is set `SUNW_PKG_ALLZONES=false`.

The following steps are performed by the `pkgadd` utility:

- If the package contents affect any area of the global zone that is shared with any non-global zone, or if the package is set `SUNW_PKG_ALLZONES=true`, then `pkgadd` fails. The error message states that the package must be added to the global zone and to all non-global zones.
- Package dependencies are checked on the global zone only. If required packages are not installed, then the dependency check fails. The system notifies the global administrator, who is prompted whether to continue.
- The package is added to the global zone.
- The package database on the global zone is updated.

- The package information on the global zone is annotated to indicate that this package is installed on the global zone only. If a non-global zone is installed in the future, this package will not be installed.

Adding a Package to all Non-Global Zones Only

To add a package to all non-global zones only, execute the `pkgadd` utility in the global zone. As the global administrator, run `pkgadd` with the `-Z` option only. Do not use the `-G` option.

To add the package to all of the non-global zones, the package must already have been installed on the global zone by using the `pkgadd` utility with the `-G` option.

The following steps are performed by the `pkgadd` utility:

- If the package is not installed on the global zone, then `pkgadd` fails. The global administrator is notified that the package must be added to the global zone and to all non-global zones.
- Package dependencies are checked on all non-global zones. If required packages are not installed, then the dependency check fails. The system notifies the global administrator, who is prompted whether to continue.
- For each non-global zone, the following actions occur:
 - The package is added to the zone.
 - The package database on the zone is updated.
- The package information on the global zone is annotated to indicate that this package is not installed in the global zone only. This package will be installed in any non-global zone that is installed in the future.

Using `pkgadd` in a Non-Global Zone

To add a package in a specified non-global zone, execute the `pkgadd` utility, without options, as the zone administrator. The following conditions apply:

- The `pkgadd` utility can only add packages in the non-global zone in which the utility is used.
- The package cannot affect any area of the zone that is shared from the global zone.
- The package must be set `SUNW_PKG_ALLZONES=false`.
- Neither the `-G` option nor the `-Z` option can be used. If either of these options is used, `pkgadd` outputs an error message and the attempted operation fails.

The following steps are performed by the `pkgadd` utility:

- Package dependencies are checked on the non-global zone's package database before the package is added. If required packages are not installed, then the dependency check fails. The system notifies the global administrator, who is prompted whether to continue. The check fails if either of the following conditions are true.
 - Any component of the package affects any area of the zone that is shared from the global zone.
 - The package is set `SUNW_PKG_ALLZONES=true`.
- The package is added to the zone.
- The package database on the zone is updated.

About Removing Packages in Zones

The `pkgrm` utility described in the `pkgrm(1M)` man page supports removing packages in a zones environment.

Using `pkgrm` in the Global Zone

The following `pkgrm` utility options are used from the global zone:

- G Remove package(s) from the global zone only.
- Z Remove package(s) from the non-global zones only

When the `pkgrm` utility is used in the global zone, the following actions apply.

- `pkgrm` can remove a package from the global zone and from all non-global zones, from all non-global zones only, or from the global zone only when the package is only installed in the global zone.
- `pkgrm` cannot remove a package from the global zone if the package is also installed in a non-global zone, or remove a package from any subset of the non-global zones.
- If the `-G` option is used, `pkgrm` removes the specified package from the global zone only.
- If the `-Z` option is used, `pkgrm` removes the specified package from all non-global zones only. The package is marked as installed in the global zone only. The package is not installed when any non-global zone is installed.
- If neither the `-G` option nor the `-Z` option is used, `pkgrm` removes the specified package from all zones, including the global zone. This is the default action.

Note that a package can only be removed from a non-global zone by a zone administrator working in that zone if the following are true:

- The package does not affect any area on the non-global zone that is shared from the global zone.
- The package is set `SUNW_PKG_ALLZONES=false`.

Removing a Package From the Global Zone and From all Non-Global Zones

To remove a package from the global zone and from all non-global zones, execute the `pkgrm` utility in the global zone. As the global administrator, run `pkgrm` without the `-G` option or the `-Z` option.

A package can be removed from the global zone and from all non-global zones without regard to the area affected by the package.

The following steps are performed by the `pkgrm` utility:

- Package dependencies are checked on the global zone and on all non-global zones. If the dependency check fails, then `pkgrm` fails. The system notifies the global administrator, who is prompted whether to continue.
- The package is removed from each non-global zone.
- The package database on each non-global zone is updated.
- The package is removed from the global zone.
- The package database on the global zone is updated.

Removing a Package From all Non-Global Zones

To remove a package from all non-global zones only, execute `pkgrm` in the global zone. As the global administrator, run `pkgrm` with the `-Z` option only.

To remove the package from all of the non-global zones, the package must already be installed on the global zone.

The following steps are performed by the `pkgrm` utility:

- If the package is not installed on the global zone, then `pkgrm` fails.
- The package dependencies are checked on all non-global zones. If a dependency check fails, then `pkgrm` fails and the global administrator is notified.
- For each non-global zone, the following actions occur:
 - The package is removed from the zone.
 - The package database on the zone is updated.

- The package information on the global zone is annotated to indicate that this package is installed in the global zone only, and that this package must not be installed on a non-global zone whenever a non-global zone is installed.

Using `pkgrm` in a Non-Global Zone

As the zone administrator, use the `pkgrm` utility in a non-global zone to remove a package. The following limitations apply:

- `pkgrm` can only remove packages from the non-global zone.
- Neither the `-G` or the `-Z` options can be used in a non-global zone. If either of these options is used, `pkgrm` outputs an error message and the attempted operation fails.
- The package cannot affect any area of the zone that is shared from the global zone.
- The package must be set `SUNW_PKG_ALLZONES=false`.

The following steps are performed by the `pkgrm` utility:

- Dependencies are checked on the non-global zone's package database. If the dependency check fails, then `pkgrm` fails and the zone administrator is notified. The check fails if either of the following conditions are true.
 - Any component of the package affects any area of the zone that is shared from the global zone.
 - The package is set `SUNW_PKG_ALLZONES=true`.
- The package is removed from the zone.
- The package database on the zone is updated.

Package Parameter Information

`SUNW_PKG_ALLZONES` Package Parameter

The optional `SUNW_PKG_ALLZONES` package parameter describes the zone scope of a package. This parameter defines the following:

- Whether a package is required to be installed on all zones
- Whether a package is required to be identical in all zones

The `SUNW_PKG_ALLZONES` package parameter has two permissible values. These values are `true` and `false`. The default value is `false`. If this parameter is either not set or set to a value other than `true` or `false`, the value `false` is used.

The `SUNW_PKG_ALLZONES` package parameter values are described in the following table.

TABLE 23-1 `SUNW_PKG_ALLZONES` Package Parameter Values

Value	Description
false	<p>This package can be installed from the global zone to the global zone only, or to the global zone and to all non-global zones. The package can also be installed from any non-global zone to the same non-global zone.</p> <ul style="list-style-type: none"> ■ The global administrator can install the package on the global zone only. ■ The global administrator can install the package on the global zone and on all non-global zones. ■ The zone administrator can install the package on a non-global zone. <p>If removed from the global zone, the package is not removed from other zones. The package can be removed from individual non-global zones.</p> <ul style="list-style-type: none"> ■ The package is not required to be installed on the global zone. ■ The package is not required to be installed on any non-global zone. ■ The package is not required to be identical across all zones. Different versions of the package can exist on individual zones. ■ The package delivers software that is not implicitly shared across all zones. This means that the package is not operating system-specific. Most application-level software is in this category. Examples include the StarOffice™ product or a web server.

TABLE 23-1 SUNW_PKG_ALLZONES Package Parameter Values (Continued)

Value	Description
true	<p>If installed on the global zone, this package must also be installed on all non-global zones. If removed from the global zone, the package must also be removed from all non-global zones.</p> <ul style="list-style-type: none">■ If the package is installed, it must be installed on the global zone. The package is then automatically installed on all non-global zones.■ The version of the package must be identical on all zones.■ The package delivers software that is implicitly shared across all zones. The package is dependent on the versions of software that are implicitly shared across all zones. The package should be visible in all non-global zones. Examples include kernel modules. These packages allow the non-global zone to resolve dependencies on packages that are installed in the global zone by requiring that the entire package be installed on all non-global zones.■ Only the global administrator can install the package. A zone administrator cannot install the package on a non-global zone.

SUNW_PKG_HOLLOW Package Parameter

The SUNW_PKG_HOLLOW package parameter defines whether a package should be visible in any non-global zone if that package is required to be installed and be identical in all zones.

The SUNW_PKG_HOLLOW package parameter has two permissible values, `true` or `false`.

- If SUNW_PKG_HOLLOW is either not set or set to a value other than `true` or `false`, the value `false` is used.
- If SUNW_PKG_ALLZONES is set to `false`, the SUNW_PKG_HOLLOW parameter is ignored.
- If SUNW_PKG_ALLZONES is set to `false`, then SUNW_PKG_HOLLOW cannot be set to `true`.

The SUNW_PKG_HOLLOW package parameter values are described in the following table.

TABLE 23-2 SUNW_PKG_HOLLOW Package Parameter Values

Value	Description
false	<p>This is not a “hollow” package:</p> <ul style="list-style-type: none"> ■ If installed on the global zone, the package content and installation information are required on all non-global zones. ■ The package delivers software that should be visible in all non-global zones. An example is the package that delivers the <code>truss</code> command. ■ Other than the restrictions for the current setting of the <code>SUNW_PKG_ALLZONES</code> package parameter, no additional restrictions are defined.
true	<p>This is a “hollow” package:</p> <ul style="list-style-type: none"> ■ The package content is not delivered on any non-global zone. However, the package installation information is required on all non-global zones. ■ The package delivers software that should not be visible in all non-global zones. Examples include kernel drivers and system configuration files that work only in the global zone. This setting allows the non-global zone to resolve dependencies on packages that are installed only on the global zone without actually installing the package data. ■ This package setting includes all of the restrictions defined for setting <code>SUNW_PKG_ALLZONES</code> to <code>true</code>. ■ In the global zone, the package is recognized as having been installed, and all components of the package are installed. Directories are created, files are installed, and class action and other scripts are run as appropriate when the package is installed. ■ In a non-global zone, the package is recognized as having been installed, but no components of the package are installed. No directories are created, no files are installed, and no class action or other install scripts are run when the package is installed. ■ When removed from the global zone, the package is recognized as having been completely installed. Appropriate directories and files are removed, and class action or other install scripts are run when the package is removed. ■ When removed from a non-global zone, the package is recognized as not having been completely installed. No directories are removed, no files are removed, and no class action or other install scripts are run when the package is removed. ■ The package is recognized as being installed in all zones for purposes of dependency checking by other packages that rely on this package being installed.

Package Information Query

The `pkginfo` utility described in the `pkginfo(1)` man page supports querying the software package database in a zones environment. For information about the database, see “[Product Database](#)” on page 281.

The following options are used with the `pkginfo` utility:

- Z Query the software package database in all non-global zones.
- z *zonename* Query the software package database in the specified non-global zone.

Using `pkginfo` in the Global Zone

When the `pkginfo` utility is used in the global zone, the following actions apply:

- `pkginfo` can query the software package database in in the global zone only, in a specified non-global zone only, or in all non-global zones only.
- `pkginfo` is not able to query the software package database in the global zone and in all non-global zones, or in the global zone and in a subset of all non-global zones.
- If the `-z zonename` option is used, `pkginfo` queries the software package database in the specified non-global zone only.
- If the `-Z` option is used, `pkginfo` queries the software package database in all non-global zones.
- If neither the `-z` option nor the `-Z` option is used, `pkginfo` queries the software package database in the global zone only.

Using `pkginfo` in a Non-Global Zone

When the `pkginfo` utility is used in a non-global zone, the following actions apply.

- `pkginfo` can query only the software package database in the non-global zone.
- Neither the `-z` option nor the `-Z` option can be used in a non-global zone. If either option is used, `pkginfo` outputs an error message and the attempted operation fails.

About Adding Patches in Zones

In general, a patch consists of the following components:

- Patch information:
 - Identification, which is the patch version and patch ID
 - Applicability, which is the operating system type, operating system version, and architecture
 - Dependencies, such as requires and obsoletes
 - Properties, such as requires a reboot afterwards
- One or more packages to patch, where each package contains:
 - The version of the package to which the patches can be applied
 - Patch information, such as ID, obsoletes, and requires
 - One or more components of the package to be patched

When the `patchadd` command is used to apply a patch, the patch information is used to determine whether the patch is applicable to the currently running system. If determined to be not applicable, the patch is not applied. Patch dependencies are also checked against all of the zones on the system. If any required dependencies are not met, the patch is not applied. This could include the case in which a later version of the patch is already installed.

Each package contained in the patch is checked. If the package is not installed on any zone, then the package is bypassed and not patched.

If all dependencies are satisfied, all packages in the patch that are installed on any zone are used to patch the system. The package and patch databases are also updated.

Applying Patches in a Zones Environment

All patches applied at the global zone level are applied across all zones. When a non-global zone is installed, it is at the same patch level as the global zone. When the global zone is patched, all non-global zones are similarly patched. This action maintains the same patch level across all zones.

The `patchadd` system utility described in the `patchadd(1M)` man page is used to add patches in a zones environment.

Using patchadd in the Global Zone

To add a patch to the global zone and to all non-global zones, run `patchadd` as the global administrator in the global zone.

When `patchadd` is used in the global zone, the following conditions apply:

- The `patchadd` utility is able to add the patch(es) to the global zone and to all non-global zones only. This is the default action.
- The `patchadd` utility cannot add the patch(es) to the global zone only or to a subset of the non-global zones.

When you add a patch to the global zone and to all non-global zones, you do not have to consider whether the patch affects areas that are shared from the global zone.

The following steps are performed by the `patchadd` utility:

- The patch is added to the global zone.
- The patch database on the global zone is updated.
- The patch is added to each non-global zone.
- The patch database on each non-global zone is updated.

Using patchadd in a Non-Global Zone

When used in a non-global zone by the zone administrator, `patchadd` can only be used to add patches to that zone. A patch can be added to a non-global zone in the following cases:

- The patch does not affect any area of the zone that is shared from the global zone.
- All packages in the patch are set `SUNW_PKG_ALLZONES=false`.

The following steps are performed by the `patchadd` utility:

- The patch is added to the zone.
- The patch database on the zone is updated.

Removing Patches in a Zones Environment

The `patchrm` system utility described in the `patchrm(1M)` man page is used to remove patches in a zones environment.

Using `patchrm` in the Global Zone

As the global administrator, you can use the `patchrm` utility in the global zone to remove patches. The `patchrm` utility cannot remove patches from the global zone only or from a subset of the non-global zones.

Using `patchrm` in a Non-Global Zone

As the zone administrator, you can use the `patchrm` utility in a non-global zone to remove patches from that non-global zone only. Patches cannot affect areas that are shared.

PatchPro Support

PatchPro can be used in the global zone and in any non-global zone. If run in the global zone, PatchPro uses the existing patch database and patch tools to patch the global and all non-global zones for all software that is installed on the global zone. No software installed in a non-global zone that is not also installed in the global zone will be taken into account.

A zone administrator can run PatchPro in a non-global zone to patch the software installed in the non-global zone.

Product Database

Each zone's respective package, patch, and product registry database completely describes all installed software that is available on the zone. All dependency checking for installing additional software or patches is performed without accessing any other zone's database, unless a package or patch is being installed or removed on the global zone and on one or more non-global zones. In this case, the appropriate non-global zone database(s) must be accessed.

For more information about the database, see the `pkgadm(1M)` man page.

Adding and Removing Packages and Patches in a Solaris Zones Environment (Tasks)

This chapter describes how to add and remove packages and patches in a Solaris zones environment. Other tasks associated with managing packages and patches, such as checking package parameter settings and obtaining package information, are also addressed. For an overview of patching and packaging concepts in a zones environment, see [Chapter 23](#).

Adding and Removing Packages and Patches in a Zones Environment (Task Map)

Task	Description	For Instructions
Add a package.	Add a package in a zones environment.	“Adding a Package in a Zones Environment” on page 284
Check package information.	Check package information in a zones environment.	“Checking Package Information in a Zones Environment” on page 286
Remove a package.	Remove a package in a zones environment.	“Removing a Package in a Zones Environment” on page 287
Apply a patch.	Apply a patch in a zones environment.	“Applying a Patch in a Zones Environment” on page 289

Task	Description	For Instructions
Remove a patch.	Remove a patch in a zones environment.	“Removing a Patch in a Zones Environment” on page 290
(Optional) Check the package parameter settings.	When adding or removing packages, verify that the settings of the package parameters support the action you want to perform.	“Checking Package Parameter Settings in a Zones Environment” on page 291

Adding a Package in a Zones Environment

You can use the `pkgadd` system utility described in the `pkgadd(1M)` man page to perform the following tasks:

- Add a package to the global zone only
- Add a package to both the global zone and all non-global zones
- Add a package that is already installed in the global zone to the non-global zones
- Add a package to a specified non-global zone only

The `SUNW_PKG_ALLZONES` and `SUNW_PKG_HOLLOW` package parameter settings must match the correct value, either `true` or `false`, to add packages in a zone environment. Otherwise, the desired result will not be achieved. For more information about the effect of these package parameter settings, see [“About Packages and Zones” on page 267](#). For more information about how to check these package parameter settings, see [“Checking Package Parameter Settings in a Zones Environment” on page 291](#).

▼ How to Add a Package to the Global Zone Only

To add a package to the global zone only, the `SUNW_PKG_ALLZONES` package parameter must be set to `false`.

You must be the global administrator in the global zone to perform this procedure.

1. **Become superuser, or assume the Primary Administrator role.**
To create the role and assign the role to a user, see [“Using the Solaris Management Tools With RBAC \(Task Map\)” in *System Administration Guide: Basic Administration*](#).
2. **While in the global zone, run the `pkgadd -d` command followed by the location of the package, the `-G` option, and then the package name.**

- If installing the package from a CD-ROM, type:

```
global# pkgadd -d /cdrom/cdrom0/directory -G package_name
```

- If installing the package from a directory to which it has been copied, type:

```
global# pkgadd -d disk1/image -G package_name
```

where *disk1* is the location where the package was copied.

▼ How to Add a Package to the Global Zone and All Non-Global Zones

Do not use `pkgadd` options `-G` and `-Z` in this procedure.

You must be the global administrator in the global zone to perform this procedure.

1. Become superuser, or assume the Primary Administrator role.

To create the role and assign the role to a user, see “Using the Solaris Management Tools With RBAC (Task Map)” in *System Administration Guide: Basic Administration*.

2. While in the global zone, run the `pkgadd -d` command followed by the location of the package and then the package name.

- If installing the package from a CD-ROM, type:

```
global# pkgadd -d /cdrom/cdrom0/directory package_name
```

- If installing the package from a directory to which it has been copied, type:

```
global# pkgadd -d disk1/image package_name
```

where *disk1* is the location where the package was copied.

▼ How to Add a Package That Is Installed in the Global Zone to All Non-Global Zones

You must be the global administrator in the global zone to perform this procedure.

1. Become superuser, or assume the Primary Administrator role.

To create the role and assign the role to a user, see “Using the Solaris Management Tools With RBAC (Task Map)” in *System Administration Guide: Basic Administration*.

2. While in the global zone, run the `pkgadd -d` command followed by the location of the package, the `-Z` option, and then the package name.

- If installing the package from a CD-ROM, type:

```
global# pkgadd -d /cdrom/cdrom0/directory -Z package_name
```

- If installing the package from a directory to which it has been copied, type:

```
global# pkgadd -d disk1/image -Z package_name
```

where *disk1* is the location where the package was copied.

▼ How to Add a Package To a Specified Non-Global Zone Only

To add a package to a specified non-global zone only, the `SUNW_PKG_ALLZONES` package parameter must be set to `false`. Do not use the `pkgadd` options `-G` and `-Z` in this procedure or the operation fails.

You must be the zone administrator in the non-global zone to perform this procedure.

1. **Log in to the non-global zone as the zone administrator.**
2. **While in the non-global zone, my-zone in this procedure, run the `pkgadd -d` command followed by the location of the package and then the package name.**

- If installing the package from a CD-ROM, type:

```
my-zone# pkgadd -d /cdrom/cdrom0/directory package_name
```

- If installing the package from a directory to which it has been copied, type:

```
my-zone# pkgadd -d disk1/image package_name
```

where *disk1* is the location where the package was copied.

Checking Package Information in a Zones Environment

You can query the software package database for the global zone and non-global zones by using the `pkginfo` command with either no option, `-Z`, or `-z zonename`. See the `pkginfo(1)` man page for more information about this command.

▼ How to Check Package Information in the Global Zone Only

- **To check the software package database for the global zone only, use `pkginfo` followed by the package name.**

```
% pkginfo package_name
```

▼ How to Check Package Information in All Non-Global Zones Only

- To check the software package database for all non-global zones, use `pkginfo -z` followed by the package name.

```
% pkginfo -z package_name
```

▼ How to Check Package Information in a Specified Non-Global Zone Only

- To check the software package database for a specific non-global zone, use `pkginfo -z` followed by the name of the zone and then the package name.

```
% pkginfo -z zonename package_name
```

Removing a Package in a Zones Environment

You can use the `pkgrm` system utility described in the `pkgrm(1M)` man page to perform the following tasks:

- Remove a package from the global zone and all non-global zones
- Remove a package from all non-global zones only
- Remove a package from a specified non-global zone only

The `SUNW_PKG_ALLZONES` and `SUNW_PKG_HOLLOW` package parameter settings must match the correct value, either `true` or `false`, to remove packages in a zones environment. Otherwise, the desired result will not be achieved. For more information about the effect of these package parameter settings, see [“About Packages and Zones” on page 267](#). For more information about how to check these package parameter settings, see [“Checking Package Parameter Settings in a Zones Environment” on page 291](#).

▼ How to Remove a Package From the Global Zone and All Non-Global Zones

To remove a package from the global zone and all non-global zones, the `SUNW_PKG_ALLZONES` package parameter must be set to `true`. Do not use `pkgrm` option `-z` in this procedure.

You must be the global administrator in the global zone to perform this procedure.

1. **Become superuser, or assume the Primary Administrator role.**

To create the role and assign the role to a user, see “Using the Solaris Management Tools With RBAC (Task Map)” in *System Administration Guide: Basic Administration*.

2. **While in the global zone, run the `pkgrm` command followed by the package name.**

```
global# pkgrm package_name
```

▼ How to Remove a Package From All Non-Global Zones

To remove a package from all non-global zones, the `SUNW_PKG_ALLZONES` package parameter must be set to `false`.

You must be the global administrator in the global zone to perform this procedure.

Note – The package must already be installed in the global zone as well or `pkgrm` fails. Likewise, if the package dependency check fails, `pkgrm` fails.

1. **Become superuser, or assume the Primary Administrator role.**

To create the role and assign the role to a user, see “Using the Solaris Management Tools With RBAC (Task Map)” in *System Administration Guide: Basic Administration*.

2. **While in the global zone, run the `pkgrm -Z` command followed by the package name.**

```
global# pkgrm -Z package_name
```

▼ How to Remove a Package From a Specified Non-Global Zone Only

To remove a package from a specified non-global zone only, the `SUNW_PKG_ALLZONES` package parameter must be set to `false`. Do not use `pkgrm` option `-Z` in this procedure.

You must be the zone administrator in the non-global zone to perform this procedure.

1. **Log in to the non-global zone as the zone administrator.**

2. **While in the non-global zone, `my-zone` in this procedure, run the `pkgrm` command followed by the package name.**

```
my-zone# pkgrm package_name
```

Applying a Patch in a Zones Environment

You can use the `patchadd` system utility described in the `patchadd(1M)` man page to perform the following tasks:

- Apply a patch to the global zone only
- Apply a patch to the global zone and all non-global zones
- Apply a patch to specified non-global zone only

▼ How to Apply a Patch to the Global Zone Only

You must be the global administrator in the global zone to perform this procedure. Do not use `pkgrm` option `-Z` in this procedure.

1. **Become superuser, or assume the Primary Administrator role.**

To create the role and assign the role to a user, see “Using the Solaris Management Tools With RBAC (Task Map)” in *System Administration Guide: Basic Administration*.

2. **While in the global zone, run the `patchadd` command followed by the patch ID.**

```
global# patchadd patch_id
```

▼ How to Apply a Patch to the Global Zone and All Non-Global Zones

You must be the global administrator in the global zone to perform this procedure.

1. **Become superuser, or assume the Primary Administrator role.**

To create the role and assign the role to a user, see “Using the Solaris Management Tools With RBAC (Task Map)” in *System Administration Guide: Basic Administration*.

2. **While in the global zone, run the `patchadd -Z` command followed by the patch ID.**

```
global# patchadd -Z patch_id
```

▼ How to Apply a Patch to a Specified Non-Global Zone Only

To apply a patch to a specified non-global zone only, the `SUNW_PKG_ALLZONES` package parameter for all packages in the patch set must be set to `false`.

You must be the zone administrator in the non-global zone to perform this procedure.

1. **Log in to the non-global zone as the zone administrator.**
2. **While in the non-global zone, my-zone in this procedure, run the `patchadd` command followed by the patch ID.**

```
my-zone# patchadd patch_id
```

Removing a Patch in a Zones Environment

You can use the `patchrm` system utility described in the `patchrm(1M)` man page to perform the following task:

- Remove a patch from a specified non-global zone only

▼ How to Remove a Patch From a Specified Non-Global Zone Only

To remove a patch from a specified non-global zone only, the `SUNW_PKG_ALLZONES` package parameter for all packages in the patch set must be set to `false`.

You must be the zone administrator in the non-global zone to perform this procedure.

1. **Log in to the non-global zone as the zone administrator.**
2. **While in the non-global zone, my-zone in this procedure, run the `patchrm` command followed by the patch ID.**

```
my-zone# patchrm patch_id
```

Checking Package Parameter Settings in a Zones Environment

Before you add or remove a software package, you can use the `pkgparam` command to check package parameter settings. This step is optional. This check also can be done when troubleshooting why a package is not added or removed as expected. For information about displaying package parameter values, see the `pkgparam(1)` man page.

▼ (Optional) How to Check the Setting of a Package Already Installed on the System

- To check the package parameter setting of a package that is already installed in a global or non-global zone, use `pkgparam` followed by the package name and the name of the parameter.

```
% pkgparam package_name SUNW_PKG_ALLZONES
true
% pkgparam package_name SUNW_PKG_HOLLOW
false
```

▼ (Optional) How to Check the Setting of a Package in Software on a CD-ROM

- To check the package parameter settings of an uninstalled package in software located on a CD-ROM, use `pkgparam -d` with the path of the CD-ROM followed by the package name and the name of the parameter.

```
% pkgparam -d /cdrom/cdrom0/directory package_name SUNW_PKG_ALLZONES
true
% pkgparam -d /cdrom/cdrom0/directory package_name SUNW_PKG_HOLLOW
false
```


Solaris Zones Administration (Overview)

This chapter covers these general zone administration topics:

- “Global Zone Visibility and Access” on page 293
- “Process ID Visibility in Zones” on page 294
- “System Observability in Zones” on page 294
- “Non-Global Zone Node Name” on page 295
- “File Systems and Non-Global Zones” on page 295
- “Networking in Non-Global Zones” on page 301
- “Device Use in Non-Global Zones” on page 303
- “Running Applications in Non-Global Zones” on page 305
- “Resource Controls Used in Non-Global Zones” on page 305
- “Fair Share Scheduler in a Zones Environment” on page 306
- “Extended Accounting in a Zones Environment” on page 307
- “Privileges in a Non-Global Zone” on page 307
- “Using IP Security Architecture in Zones” on page 308
- “Using Solaris Auditing in Zones” on page 308
- “Core Files in Zones” on page 310
- “Commands Used in the Solaris Zones Environment” on page 310

Global Zone Visibility and Access

The global zone acts as both the default zone for the system and as a zone for system-wide administrative control. There are administrative issues associated with this dual role. Since applications within the zone have access to processes and other system objects in other zones, the effect of administrative actions can be wider than expected. For example, service shutdown scripts often use `pkill` to signal processes of a given name to exit. When such a script is run from the global zone, all such processes in the system will be signaled, regardless of zone.

The system-wide scope is often needed. For example, to monitor system-wide resource usage, you must view process statistics for the whole system. A view of just global zone activity would miss relevant information from other zones in the system that might be sharing some or all of the system resources. Such a view is particularly important when system resources such as CPU are not strictly partitioned using resource management facilities.

Thus, processes in the global zone can observe processes and other objects in non-global zones. This allows such processes to have system-wide observability. The ability to control or send signals to processes in other zones is restricted by the privilege `PRIV_PROC_ZONE`. The privilege is similar to `PRIV_PROC_OWNER` because the privilege allows processes to override the restrictions placed on unprivileged processes. In this case, the restriction is that unprivileged processes in the global zone cannot signal or control processes in other zones. This is true even when the user IDs of the processes match or the acting process has the `PRIV_PROC_OWNER` privilege. The `PRIV_PROC_ZONE` privilege can be removed from otherwise privileged processes to restrict actions to the global zone.

For information about matching processes by using a `zoneidlist`, see the `pgrep(1)` and `pkill(1)` man pages.

Process ID Visibility in Zones

Only processes in the same zone will be visible through system call interfaces that take process IDs, such as the `kill` and `priocntl` commands. For information, see the `kill(1)` and the `priocntl(1)` man pages.

System Observability in Zones

The `ps` command has the following modifications:

- The `-o` option is used to specify output format. This option allows you to print the zone ID of a process or the name of the zone in which the process is running.
- The `-z zonelist` option is used to list only processes in the specified zones. Zones can be specified either by zone name or by zone ID. This option is only useful when the command is executed in the global zone.
- The `-Z` option is used to print the name of the zone associated with the process. The name is printed under the column heading `ZONE`.

For more information, see the `ps(1)` man page.

A `-z zonename` option has been added to the following Solaris utilities. You can use this option to filter the information to include only the zone or zones specified.

- `ipcs` (see the `ipcs(1)` man page)
- `pgrep` (see the `pgrep(1)` man page)
- `ptree` (see the `proc(1)` man page)
- `prstat` (see the `prstat(1M)` man page)

See [Table 25-2](#) for the full list of changes made to commands.

Non-Global Zone Node Name

The node name in `/etc/nodename` returned by `uname -n` can be set by the zone administrator. The node name must be unique.

File Systems and Non-Global Zones

This section provides information about file system issues in a zones environment. Each zone has its own section of the file system hierarchy, rooted at a directory known as the zone `root`. Processes in the zone can access only files in the part of the hierarchy that is located under the zone `root`. The `chroot` utility can be used in a zone, but only to restrict the process to a root path within the zone. For more information about `chroot`, see `chroot(1M)`.

The `-o nosuid` Option

The `-o nosuid` option to the `mount` utility has the following functionality:

- Processes from a `setuid` binary located on a file system that is mounted using the `nosetuid` option do not run with the privileges of the `setuid` binary. The processes run with the privileges of the user that executes the binary.
For example, if a user executes a `setuid` binary that is owned by `root`, the processes run with the privileges of the user.
- Opening device-special entries in the file system is not allowed. This behavior is equivalent to specifying the `nodevices` option.

This file system-specific option is available to all Solaris file systems that have mount utilities. These file systems are listed in [“Mounting File Systems in Zones”](#) on page 296. Mounting capabilities are also described. For more information about the `-o nosuid` option, see [“Accessing Network File Systems \(Reference\)”](#) in *System Administration Guide: Network Services*.

Mounting File Systems in Zones

Options for mounting file systems in non-global zones are described in the following table. Procedures for these mounting alternatives are provided in [“Configuring, Verifying, and Committing a Zone”](#) on page 226 and [“Mounting File Systems in Running Non-Global Zones”](#) on page 319.

Note – When file systems are mounted from within a zone, the `nodevices` option applies. For example, if a zone is granted access to a block device (`/dev/dsk/c0t0d0s7`) and a raw device (`/dev/rdisk/c0t0d0s7`) corresponding to a UFS file system, the file system is automatically mounted `nodevices` when mounted from within a zone. This rule does not apply to mounts specified through a `zonecfg` configuration.

Note – Any file system type not listed in the table can be specified in the configuration if it has a mount binary in `/usr/lib/fstype/mount`.

File System	Mounting Options in a Non-Global Zone
AutoFS	Cannot be mounted using <code>zonecfg</code> , cannot be manually mounted from the global zone into a non-global zone. Can be mounted from within the zone.
CacheFS	Cannot be used in a non-global zone.
FDFS	Can be mounted using <code>zonecfg</code> , can be manually mounted from the global zone into a non-global zone, can be mounted from within the zone.
HSFS	Can be mounted using <code>zonecfg</code> , can be manually mounted from the global zone into a non-global zone, can be mounted from within the zone.

File System	Mounting Options in a Non-Global Zone
LOFS	Can be mounted using <code>zonecfg</code> , can be manually mounted from the global zone into a non-global zone, can be mounted from within the zone.
MNTFS	Cannot be mounted using <code>zonecfg</code> , cannot be manually mounted from the global zone into a non-global zone. Can be mounted from within the zone.
NFS	Cannot be mounted using <code>zonecfg</code> . V2, V3, and V4, which are the versions currently supported in zones, can be mounted from within the zone.
PCFS	Can be mounted using <code>zonecfg</code> , can be manually mounted from the global zone into a non-global zone, can be mounted from within the zone.
PROCFS	Cannot be mounted using <code>zonecfg</code> , cannot be manually mounted from the global zone into a non-global zone. Can be mounted from within the zone.
TMPFS	Can be mounted using <code>zonecfg</code> , can be manually mounted from the global zone into a non-global zone, can be mounted from within the zone.
UDFS	Can be mounted using <code>zonecfg</code> , can be manually mounted from the global zone into a non-global zone, can be mounted from within the zone.
UFS	Can be mounted using <code>zonecfg</code> , can be manually mounted from the global zone into a non-global zone, can be mounted from within the zone.
XMEMFS	Can be mounted using <code>zonecfg</code> , can be manually mounted from the global zone into a non-global zone, can be mounted from within the zone.

For more information, see [“How to Configure the Zone”](#) on page 226, [“Mounting File Systems in Running Non-Global Zones”](#) on page 319, and the `mount(1M)` man page.

Security Restrictions and File System Behavior

There are security restrictions on mounting certain file systems from within a zone. Other file systems exhibit special behavior when mounted in a zone. The list of modified file systems follows.

AutoFS

Autofs is a client-side service that automatically mounts the appropriate file system. When a client attempts to access a file system that is not presently mounted, the autofs file system intercepts the request and calls `automountd` to mount the requested directory. AutoFS mounts established within a zone are local to that zone. The mounts cannot be accessed from other zones, including the global zone. The mounts are removed when the zone is halted or rebooted. For more information on AutoFS, see “How Autofs Works” in *System Administration Guide: Network Services*.

Each zone runs its own copy of `automountd`. The auto maps and timeouts are controlled by the zone administrator. You cannot trigger a mount in another zone by crossing an AutoFS mount point for a non-global zone from the global zone.

Certain AutoFS mounts are created in the kernel when another mount is triggered. Such mounts cannot be removed by using the regular `umount` interface because they must be mounted or unmounted as a group. Note that this functionality is provided for zone shutdown.

MNTFS

MNTFS is a virtual file system that provides read-only access to the table of mounted file systems for the local system. The set of file systems visible by using `mnttab` from within a non-global zone is the set of file systems mounted in the zone, plus an entry for root (/). Mount points with a special device that is not accessible from within the zone, such as `/dev/rdisk/c0t0d0s0`, have their special device set to the same as the mount point. All mounts in the system are visible from the global zone’s `/etc/mnttab` table. For more information on MNTFS, see “Mounting and Unmounting File Systems (Tasks)” in *System Administration Guide: Devices and File Systems*.

NFS

NFS mounts established within a zone are local to that zone. The mounts cannot be accessed from other zones, including the global zone. The mounts are removed when the zone is halted or rebooted.

As documented in the `mount_nfs(1M)` man page, an NFS server should not attempt to mount its own file systems. Thus, a zone should not NFS mount a file system exported by the global zone. Zones cannot be NFS servers. From within a zone, NFS mounts behave as though mounted with the `nodevices` option.

The `nfsstat` command output only pertains to the zone in which the command is run. For example, if the command is run in the global zone, only information about the global zone is reported. For more information about the `nfsstat` command, see `nfsstat(1M)`.

The `zlogin` command will fail if any of its open files or any portion of its address space reside on NFS. For more information, see [“zlogin Command” on page 253](#).

PROCFS

The `/proc` file system, or PROCFS, provides process visibility and access restrictions as well as information about the zone association of processes. Only processes in the same zone are visible through `/proc`.

Processes in the global zone can observe processes and other objects in non-global zones. This allows such processes to have system-wide observability.

From within a zone, `procfs` mounts behave as though mounted with the `nodevices` option. For more information about `procfs`, see the `proc(4)` man page.

LOFS

The scope of what can be mounted through LOFS is limited to the portion of the file system that is visible to the zone. Hence, there are no restrictions on LOFS mounts in a zone.

UFS, UDFS, HSFs, PCFS, and other storage-based file systems

When using the `zonecfg` command to configure storage-based file systems that have an `fsck` binary, such as UFS, the zone administrator must specify a `raw` parameter. The parameter indicates the raw (character) device, such as `/dev/rdisk/c0t0d0s7`. `zoneadmd` automatically runs the `fsck` command in non-interactive check-only mode (`fsck -m`) on this device before it mounts the file system. If the `fsck` fails, `zoneadmd` cannot bring the zone to the ready state. The path specified by `raw` cannot be a relative path.

It is an error to specify a device to `fsck` for a file system that does not provide an `fsck` binary in `/usr/lib/fstype/fsck`. It is also an error if you do not specify a device to `fsck` if an `fsck` binary exists for that file system.

For more information, see [“The zoneadmd Daemon” on page 239](#) and the `fsck(1M)` man page.

Non-Global Zones as NFS Clients

Zones can be NFS clients. Version 2, version 3, and version 4 protocols are supported. For information on these NFS versions, see [“Features of the NFS Service” in *System Administration Guide: Network Services*](#).

The default version is NFS version 4. You can enable other NFS versions on a client by using one of the following methods:

- You can edit `/etc/default/nfs` to set `NFS_CLIENT_VERSMAX=number` so that the zone uses the specified version by default. See [“Task Map for NFS Services” in *System Administration Guide: Network Services*](#). Use the procedure [“How to Select Different Versions of NFS on a Client by Modifying the `/etc/default/nfs` File”](#)

from the task map.

- You can manually create a version mount. This method overrides the contents of `/etc/default/nfs`. See “Task Map for NFS Services” in *System Administration Guide: Network Services*. Use the procedure “How to Use the Command Line to Select Different Versions of NFS on a Client” from the task map.

Use of `mknod` Prohibited in a Zone

Note that you cannot use the `mknod` command documented in the `mknod(1M)` man page to make a special file in a non-global zone.

Traversing File Systems

A zone’s file system namespace is a subset of the namespace accessible from the global zone. Unprivileged processes in the global zone are prevented from traversing a non-global zone’s file system hierarchy through the following means:

- Specifying that the zone root’s parent directory is owned, readable, writable, and executable by root only
- Restricting access to directories exported by `/proc`

Note that attempting to access AutoFS nodes mounted for another zone will fail. The global administrator must not have auto maps that descend into other zones.

Restriction on Accessing A Non-Global Zone From the Global Zone

After a non-global zone is installed, the zone must never be accessed directly from the global zone by any commands other than system backup utilities. Moreover, a non-global zone can no longer be considered secure after it has been exposed to an unknown environment. An example would be a zone placed on a publicly accessible network, where it would be possible for the zone to be compromised and the contents of its file systems altered. If there is any possibility that compromise has occurred, the global administrator should treat the zone as untrusted.

Any command that accepts an alternative root by using the `-R` option (or the equivalent) must not be used when the following are true:

- The command is run in the global zone.
- The alternative root is a non-global zone.

An example is the `-R root_path` option to the `pkgadd` utility run from the global zone with a non-global zone root path.

Networking in Non-Global Zones

In a zones environment, the zones can communicate with each other over the network. The zones all have separate bindings, or connections, and the zones can all run their own server daemons. These daemons can listen on the same port numbers without any conflict. The IP stack resolves conflicts by considering the IP addresses for incoming connections. The IP addresses identify the zone.

Zone Partitioning

The IP stack in a system supporting zones implements the separation of network traffic between zones. Applications that receive IP traffic can only receive traffic sent to the same zone.

Each logical interface on the system belongs to a specific zone, the global zone by default. Logical network interfaces assigned to zones through the `zonecfg` utility are used to communicate over the network. Each stream and connection belongs to the zone of the process that opened it.

Bindings between upper-layer streams and logical interfaces are restricted. A stream can 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.

Each zone has its own set of binds. Each zone can be running the same application listening on the same port number without binds failing because the address is already in use. Each zone can run its own version of the following services:

- Internet services daemon with a full configuration file (see the `inetd(1M)` man page)
- `sendmail` (see the `sendmail(1M)` man page)
- `apache` (see the `apache(1M)` man page)

Zones other than the global zone have restricted access to the network. The standard TCP and UDP socket interfaces are available, but `SOCK_RAW` socket interfaces are restricted to Internet Control Message Protocol (ICMP). ICMP is necessary for detecting and reporting network error conditions or using the `ping` command.

Network Interfaces

Each non-global zone that requires network connectivity has one or more dedicated IP addresses. These addresses are associated with logical network interfaces that can be placed in a zone by using the `ifconfig` command. Zone interfaces configured by

`zonecfg` will automatically be plumbed and placed in the zone when it is booted. The `ifconfig` command can be used to add or remove logical interfaces when the zone is running. Only the global administrator can modify the interface configuration and the network routes.

Within a non-global zone, only that zone's interfaces will be visible to `ifconfig`.

For more information, see the `ifconfig(1M)` and `if_tcp(7P)` man pages.

IP Traffic Between Zones on the Same Machine

Between two zones on the same machine, packet delivery is only allowed if there is a "matching route" for the destination and the zone in the forwarding table.

The matching information is implemented as follows:

- The source address for the packets is selected on the output interface specified by the matching route.
- By default, traffic is permitted between two zones that have addresses on the same subnet. The matching route in this case is the interface route for the subnet.
- If there is a default route for a given zone, where the gateway is on one of the zone's subnets, traffic from that zone to all other zones is allowed. The matching route in this case is the default route.
- If there is a matching route with the `RTF_REJECT` flag, packets trigger an ICMP unreachable message. If there is a matching route with the `RTF_BLACKHOLE` flag, packets are discarded. The global administrator can use the `route` command options described in the following table to create routes with these flags.

Modifier	Flag	Description
<code>-reject</code>	<code>RTF_REJECT</code>	Emit an ICMP unreachable message when matched.
<code>-blackhole</code>	<code>RTF_BLACKHOLE</code>	Silently discard packets during updates.

For more information, see the `route(1M)` man page.

IP Network Multipathing in Zones

IP network multipathing (IPMP) provides physical interface failure detection and transparent network access failover for a system with multiple interfaces on the same IP link. IPMP also provides load spreading of packets for systems with multiple interfaces.

All network configuration is done in the global zone. You can configure IPMP in the global zone, then extend the functionality to non-global zones. The functionality is extended by placing the zone's address in an IPMP group when you configure the zone. Then, if one of the interfaces in the global zone fails, the non-global zone addresses will migrate to another network interface card.

In a given non-global zone, only the interfaces associated with the zone are visible through the `ifconfig` command.

See “How to Extend IP Network Multipathing Functionality to Non-Global Zones” on page 322. The zones configuration procedure is covered in “How to Configure the Zone” on page 226. For information on IPMP features, components, and usage, see “IP Network Multipathing (Overview)” in *System Administration Guide: IP Services*.

Device Use in Non-Global Zones

The set of devices available within a zone is restricted to prevent a process in one zone from interfering with processes running in other zones. For example, a process in a zone cannot modify kernel memory or modify the contents of the root disk. Thus, by default, only certain pseudo-devices that are considered safe for use in a zone are available. Additional devices can be made available within specific zones by using the `zonectg` utility.

`/dev` and the `/devices` Namespace

The `devfs` file system described in the `devfs(7FS)` man page is used by the Solaris to manage `/devices`. Each element in this namespace represents the physical path to a hardware device, pseudo-device, or nexus device. The namespace is a reflection of the device tree. As such, the file system is populated by a hierarchy of directories and device special files.

The `/dev` file hierarchy, which is today part of the `/` (root) file system, consists of symbolic links, or logical paths, to the physical paths present in `/devices`. Applications reference the logical path to a device presented in `/dev`. The `/dev` file system is loopback-mounted into the zone using a read-only mount.

The `/dev` file hierarchy is managed by a system comprised of the components in the following list:

- `devfsadm` (see the `devfsadm(1M)` man page)
- `syseventd` (see the `syseventd(1M)` man page)
- `libdevinfo` device information library (see the `libdevinfo(3LIB)` man page)

- `devinfo` driver (see `devinfo(7D)`)
- Reconfiguration Coordination Manager (RCM) (see “Reconfiguration Coordination Manager (RCM) Script Overview” in *System Administration Guide: Devices and File Systems*)



Caution – Subsystems that rely on `/devices` path names are not able to run in non-global zones until `/dev` path names are established.

Exclusive-Use Devices

You might have devices that you want to assign to specific zones. Allowing unprivileged users to access specific devices could permit those devices to be used to cause system panic, bus resets, or other adverse effects. Before making such assignments, consider the following issues:

- Before assigning a SCSI tape device to a specific zone, consult the `sgen(7D)` man page.
- Placing a physical device into more than one zone can create a covert channel between zones. Global zone applications that use such a device risk the possibility of compromised data or data corruption by a non-global zone.

Device Driver Administration

In a non-global zone, you can use the `modinfo` command described in the `modinfo(1M)` man page to examine the list of loaded kernel modules.

Most operations concerning kernel, device, and platform management will not work inside a non-global zone because modifying platform hardware configurations violates the zone security model. These operations include the following:

- Adding and removing drivers
- Explicitly loading and unloading kernel modules
- Initiating dynamic reconfiguration (DR) operations
- Using facilities that affect the state of the physical platform

Utilities That Do Not Work or Are Modified in Non-Global Zones

Utilities That Do Not Work in Non-Global Zones

The following utilities do not work in a zone because they rely on devices that are not normally available:

- `prtconf` (see the `prtconf(1M)` man page)
- `prtdiag` (see the `prtdiag(1M)` man page)

SPARC: Utility Modified for Use in a Non-Global Zone

The `eeeprom` utility can be used in a zone to view settings. The utility cannot be used to change settings. For more information, see the `eeeprom(1M)` and `openprom(7D)` man pages.

Running Applications in Non-Global Zones

In general, all applications can run in a non-global zone. However, the following types of applications might not be suitable for this environment:

- Applications that use privileged operations that affect the system as a whole. Examples include operations that set the global system clock or lock down physical memory.
- The few applications dependent upon certain devices that do not exist in a non-global zone, such as `/dev/kmem` or `/dev/ip`.
- Applications that expect to be able to write into `/usr`, either at runtime or when being installed, patched, or upgraded. This is because `/usr` is read-only for a non-global zone by default. Sometimes the issues associated with this type of application can be mitigated without changing the application itself.

Resource Controls Used in Non-Global Zones

For additional information about using a resource management feature in a zone, also refer to the chapter that describes the feature in Part 1 of this guide.

Any of the resource controls and attributes described in the resource management chapters can be set in the non-global zone `/etc/project` file, NIS map, or LDAP directory service. The settings for a given non-global zone affect only that zone. A project running autonomously in different zones can have controls set individually in each zone. For example, Project A in the global zone can be set `project.cpu-shares=10` while Project A in a non-global zone can be set `project.cpu-shares=5`. You could have several instances of `rcapd` running, with each instance operating only on its zone.

The resource controls and attributes used in a zone to control projects, tasks, and processes within that zone are subject to the additional requirements regarding pools and the zone-wide resource controls.

A “one zone, one pool” rule applies to non-global zones. Processes in the global zone, however, can be bound by a sufficiently privileged process to any pool. The resource controller `poold` only runs in the global zone, where there is more than one pool for it to operate on. The `poolstat` utility run in a non-global zone displays only information about the pool associated with the zone. The `pooladm` command run without arguments in a non-global zone displays only information about the pool associated with the zone.

Zone-wide resource controls do not take effect when they are set in the `project` file. A zone-wide resource control is set through the `zonecfg` utility.

Fair Share Scheduler in a Zones Environment

This section describes fair share scheduler (FSS) use in a zones environment. See [“Using the Fair Share Scheduler in a Zones Environment” on page 323](#) for example procedures.

FSS Share Division in a Non-Global Zone

FSS CPU shares for a zone are hierarchical. The shares for a given non-global zone are set by the global administrator through the zone-wide resource control `zone.cpu-shares`. The `project.cpu-shares` resource control can then be defined for each project within that zone to further subdivide the shares set through the zone-wide control.

For more information on `project.cpu-shares`, see [“Available Resource Controls” on page 74](#).

Share Balance Between Zones

The global zone is given one share by default. If you have one non-global zone on your system and you give this zone two shares through `zone.cpu-shares`, that defines the proportion of CPU which the non-global zone will receive in relation to the global zone. The ratio of CPU between the two zones is 2:1.

Extended Accounting in a Zones Environment

The extended accounting subsystem collects and reports information for the entire system (including non-global zones) when run in the global zone. The global administrator can also determine resource consumption on a per-zone basis.

The extended accounting subsystem permits different accounting settings and files on a per-zone basis for process-based and task-based accounting. The `exacct` records can be tagged with the zone name `EXD PROC ZONENAME` for processes, and the zone name `EXD TASK ZONENAME` for tasks. Accounting records are written to the global zone's accounting files as well as the per-zone accounting files. The `EXD TASK HOSTNAME`, `EXD PROC HOSTNAME`, and `EXD HOSTNAME` records contain the `uname -n` value for the zone in which the process or task executed instead of the global zone's node name.

For information about IPQoS flow accounting, see “Using Flow Accounting and Statistics Gathering (Tasks)” in *System Administration Guide: IP Services*.

Privileges in a Non-Global Zone

Processes are restricted to a subset of privileges. Privilege restriction prevents a zone from performing operations that might affect other zones. The set of privileges limits the capabilities of privileged users within the zone. To display the list of privileges available within a zone, use the `ppriv` utility.

For usage examples, see “Using the `ppriv` Utility” on page 315. For more information about privileges, see the `ppriv(1)` man page and *System Administration Guide: Security Services*.

Using IP Security Architecture in Zones

The Internet Protocol Security Architecture (IPsec), which provides IP datagram protection, is described in “IP Security Architecture (Overview)” in *System Administration Guide: IP Services*. The Internet Key Exchange (IKE) protocol is used to manage the required keying material for authentication and encryption automatically.

IPsec can be used in the global zone. However, IPsec in a non-global zone cannot use IKE. Therefore, you must manage the IPsec keys and policy for the non-global zones by running the `ipseckey` and `ipseccnf` commands from the global zone. Use the source address that corresponds to the non-global zone that you are configuring.

For more information, see the `ipseccnf(1M)` and `ipseckey(1M)` man pages.

Using Solaris Auditing in Zones

Solaris auditing is described in “Solaris Auditing (Overview)” in *System Administration Guide: Security Services*. For zones considerations associated with auditing, see the following sections:

- “Planning for Solaris Auditing” in *System Administration Guide: Security Services*
- “Auditing and Zones” in *System Administration Guide: Security Services*

An audit record describes an event, such as logging in to a system or writing to a file. The record is composed of tokens, which are sets of audit data. By using the `zonename` token, you can configure Solaris auditing to identify audit events by zone. Use of the `zonename` token allows you to produce the following information:

- Audit records that are marked with the name of the zone that generated the record
- An audit log for a specific zone that the global administrator can make available to the zone administrator

Configuring Audit in the Global Zone

Solaris audit trails are configured in the global zone. Audit policy is set in the global zone and applies to processes in all zones. The audit records can be marked with the name of the zone in which the event occurred. To include zone names in audit records, you must edit the `/etc/security/audit_startup` file before you install any non-global zones. The zone name selection is case-sensitive.

To configure auditing in the global zone to include all zone audit records, add this line to the `/etc/security/audit_startup` file:

```
/usr/sbin/auditconfig -setpolicy +zonename
```

As the global administrator in the global zone, execute the `auditconfig` utility:

```
global# auditconfig -setpolicy +zonename
```

For additional information, see the `audit_startup(1M)` and `auditconfig(1M)` man pages and “Configuring Audit Files (Task Map)” in *System Administration Guide: Security Services*.

Configuring User Audit Characteristics in a Non-Global Zone

When a non-global zone is installed, the `audit_control` file and the `audit_user` file in the global zone are copied to the zone’s `/etc/security` directory. These files might require modification to reflect the zone’s audit needs.

For example, each zone can be configured to audit some users differently from others. To apply different per-user preselection criteria, both the `audit_control` and the `audit_user` files must be edited. The `audit_user` file in the non-global zone might also require revisions to reflect the user base for the zone if necessary. Because each zone can be configured differently with regard to auditing users, it is possible for the `audit_user` file to be empty.

For additional information, see the `audit_control(4)` and `audit_user(4)` man pages.

Providing Audit Records for a Specific Non-Global Zone

By including the `zonename` token as described in “Configuring Audit in the Global Zone” on page 308, Solaris audit records can be categorized by zone. Records from different zones can then be collected by using the `auditreduce` command to create logs for a specific zone.

For more information, see the `audit_startup(1M)` and `auditreduce(1M)` man pages.

Core Files in Zones

The `coreadm` command is used to specify the name and location of core files produced by abnormally terminating processes. Core file paths that include the *zonename* of the zone in which the process executed can be produced by specifying the `%z` variable. The path name is relative to a zone's root directory.

For more information, see the `coreadm(1M)` and `core(4)` man pages.

Commands Used in the Solaris Zones Environment

The commands identified in the table *Commands Used to Administer Zones* provide the primary administrative interface to the zones facility.

TABLE 25-1 Commands Used to Administer Zones

Command Reference	Description
<code>zlogin(1)</code>	Log in to a non-global zone
<code>zonename(1)</code>	Prints the name of the current zone
<code>zoneadm(1M)</code>	Administers zones on a system
<code>zonecfg(1M)</code>	Used to set up a zone configuration
<code>getzoneid(3C)</code>	Used to map between zone ID and name
<code>zones(5)</code>	Provides description of zones facility
<code>zcons(7D)</code>	Zone console device driver

The `zoneadmd` daemon is the primary process for managing the zone's virtual platform. The man page for the `zoneadmd` daemon is `zoneadmd(1M)`. The daemon does not constitute a programming interface.

The commands identified in the table *Commands Modified for use in the Zones Environment* have options that are specific to zones or present information differently in a zones environment.

TABLE 25-2 Commands Modified for use in the Zones Environment

Command Reference	Description
ipcrm(1)	Added <code>-z zone</code> option. This option is only useful when the command is executed in the global zone.
ipcs(1)	Added <code>-z zone</code> option. This option is only useful when the command is executed in the global zone.
pgrep(1)	Added <code>-z zoneidlist</code> option. This option is only useful when the command is executed in the global zone.
ppriv(1)	Added the expression <code>zone</code> for use with the <code>-l</code> option to list all privileges available in the current zone. Also use the option <code>-v</code> after <code>zone</code> to obtain verbose output.
priocntl(1)	Zone ID can be used in <code>idlist</code> and <code>-i idtype</code> to specify processes. You can use the <code>priocntl -i zoneid</code> command to move running processes into a different scheduling class in a non-global zone.
proc(1)	Added <code>-z zone</code> option to <code>ptree</code> only. This option is only useful when the command is executed in the global zone.
ps(1)	Added <code>zonename</code> and <code>zoneid</code> to list of recognized <code>format</code> names used with the <code>-o</code> option. Added <code>-z zonelist</code> to list only processes in the specified zones. Zones can be specified either by zone name or by zone ID. This option is only useful when the command is executed in the global zone. Added <code>-Z</code> to print the name of the zone associated with the process. The name is printed under an additional column header, <code>ZONE</code> .
renice(1)	Added <code>zoneid</code> to list of valid arguments used with the <code>-i</code> option.
sar(1)	If executed in a non-global zone in which the pools facility is enabled, the <code>-b</code> , <code>-c -g</code> , <code>-m</code> , <code>-p</code> , <code>-u</code> , <code>-w</code> , and <code>-y</code> options display values only for processors that are in the processor set of the pool to which the zone is bound.
auditconfig(1M)	Added <code>zonename</code> token.
auditreduce(1M)	Added <code>-z zone-name</code> option. Added ability to get an audit log of a zone.
coreadm(1M)	Added variable <code>%z</code> to identify the zone in which process executed.
df(1M)	Added <code>-Z</code> option to display mounts in all visible zones.
ifconfig(1M)	Added <code>zone</code> option for global zone use (the default), and <code>-zone zonename</code> for non-global zone use.

TABLE 25-2 Commands Modified for use in the Zones Environment (Continued)

Command Reference	Description
<code>iostat(1M)</code>	If executed in a non-global zone in which the pools facility is enabled, information is provided only for those processors that are in the processor set of the pool to which the zone is bound.
<code>mpstat(1M)</code>	If executed in a non-global zone in which the pools facility is enabled, command only displays lines for the processors that are in the processor set of the pool to which the zone is bound.
<code>poolbind(1M)</code>	Added <code>zoneid</code> list. Also see “Resource Pools Used in Zones” on page 135 for information about using zones with resource pools.
<code>prstat(1M)</code>	Added <code>-z zoneidlist</code> option. Also added <code>-Z</code> option. If executed in a non-global zone in which the pools facility is enabled, the percentage of recent CPU time used by the process is displayed only for the processors in the processor set of the pool to which the zone is bound.
<code>psrinfo(1M)</code>	If executed in a non-global zone, only information about the processors visible to the zone is displayed.
<code>traceroute(1M)</code>	Usage change. When specified from within a non-global zone, the <code>-F</code> option has no effect because the “don’t fragment” bit is always set.
<code>vmstat(1M)</code>	When executed in a non-global zone in which the pools facility is enabled, statistics are reported only for the processors in the processor set of the pool to which the zone is bound. Applies to output from the <code>-p</code> option and the <code>page</code> , <code>faults</code> , and <code>cpu</code> report fields.
<code>auditon(2)</code>	Added <code>AUDIT_ZONEID</code> to generate a zone ID token with each audit record.
<code>priocntl(2)</code>	Added <code>P_ZONEID id</code> argument.
<code>processor_info(2)</code>	If the caller is in a non-global zone and the pools facility is enabled, but the processor is not in the processor set of the pool to which the zone is bound, an error is returned.
<code>p_online(2)</code>	If the caller is in a non-global zone and the pools facility is enabled, but the processor is not in the processor set of the pool to which the zone is bound, an error is returned.
<code>pset_bind(2)</code>	Added <code>P_ZONEID</code> as <code>idtype</code> . Added <code>zone</code> to possible choices for <code>P_MYID</code> specification. Added <code>P_ZONEID</code> to valid <code>idtype</code> list in <code>EINVAL</code> error description.
<code>pset_info(2)</code>	If the caller is in a non-global zone and the pools facility is enabled, but the processor is not in the processor set of the pool to which the zone is bound, an error is returned.

TABLE 25-2 Commands Modified for use in the Zones Environment (Continued)

Command Reference	Description
pset_list(2)	If the caller is in a non-global zone and the pools facility is enabled, but the processor is not in the processor set of the pool to which the zone is bound, an error is returned.
pset_setattr(2)	If the caller is in a non-global zone and the pools facility is enabled, but the processor is not in the processor set of the pool to which the zone is bound, an error is returned.
sysinfo(2)	Changed PRIV_SYS_CONFIG to PRIV_SYS_ADMIN.
umount(2)	ENOENT is returned if file pointed to by <i>file</i> is not an absolute path.
getloadavg(3C)	If the caller is in a non-global zone and the pools facility is enabled, the behavior is equivalent to calling with a pset id of PS_MYID.
getpriority(3C)	Added zone IDs to target processes that can be specified. Added zone ID to EINVAL error description.
priv_str_to_set(3C)	Added “zone” string for the set of all privileges available within the caller’s zone.
pset_getloadavg(3C)	If the caller is in a non-global zone and the pools facility is enabled, but the processor is not in the processor set of the pool to which the zone is bound, an error is returned.
sysconf(3C)	If the caller is in a non-global zone and the pools facility enabled, sysconf(_SC_NPROCESSORS_CONF) and sysconf(_SC_NPROCESSORS_ONLN) return the number of processors in the processor set of the pool to which the zone is bound.
ucred_get(3C)	Added ucred_getzoneid() function, which returns the zone ID of the process or -1 if the zone ID is not available.
core(4)	Added n_type: NT_ZONENAME. This entry contains a string that describes the name of the zone in which the process was running.
pkginfo(4)	Now provides optional parameters and an environment variable in support of zones.
proc(4)	Added capability to obtain information on processes running in zones.
audit_syslog(5)	Added in<zone name> field that is used if the zonename audit policy is set.
privileges(5)	Added PRIV_PROC_ZONE, which allows a process to trace or send signals to processes in other zones. See zones(5).
if_tcp(7P)	Added zone ioctl() calls.

TABLE 25-2 Commands Modified for use in the Zones Environment (Continued)

Command Reference	Description
cmn_err(9F)	Added zone parameter.
ddi_cred(9F)	Added <code>crgetzoneid()</code> , which returns the zone ID from the user credential pointed to by <code>cr</code> .

Solaris Zones Administration (Tasks)

This chapter covers general administration tasks and provides usage examples.

- “Using the `ppriv` Utility” on page 315
- “Mounting File Systems in Running Non-Global Zones” on page 319
- “Using IP Network Multipathing in a Zones Environment” on page 322
- “Using the Fair Share Scheduler in a Zones Environment” on page 323
- “Using Rights Profiles in Zone Administration” on page 324

Using the `ppriv` Utility

Use the `ppriv` utility to display the zone’s privileges.

▼ How to List the Non-Global Zone’s Privilege Set

Use the `ppriv` utility with the `-l` option and the expression `zone` to list the zone’s privileges.

1. **Log into the non-global zone. This example uses a zone named `my-zone`.**
2. **At the prompt, type `ppriv -l zone` to report the set of privileges available in the zone.**

```
my-zone# ppriv -l zone
```

You will see a display similar to this:

```
file_chown
file_chown_self
file_dac_execute
file_dac_read
```

```

file_dac_search
file_dac_write
file_link_any
file_owner
file_setdac
file_setid
ipc_dac_read
ipc_dac_write
ipc_owner
net_icmpaccess
net_privaddr
proc_chroot
proc_audit
proc_exec
proc_fork
proc_owner
proc_session
proc_setid
proc_taskid
sys_acct
sys_admin
sys_mount
sys_nfs
sys_resource

```

▼ How to List a Non-Global Zone's Privilege Set With Verbose Output

Use the `ppriv` utility with the `-l` option, the expression `zone`, and the `-v` option to list the zone's privileges.

1. **Log into the non-global zone.** This example uses a zone named *my-zone*.
2. **At the prompt, type `ppriv -l zone -v` to report the set of privileges available in the zone, with a description of each privilege.**

```
my-zone# ppriv -l zone -v
```

You will see a display similar to this:

```

file_chown
    Allows a process to change a file's owner user ID.
    Allows a process to change a file's group ID to one other than
    the process' effective group ID or one of the process'
    supplemental group IDs.
file_chown_self
    Allows a process to give away its files; a process with this
    privilege will run as if {_POSIX_CHOWN_RESTRICTED} is not
    in effect.
file_dac_execute
    Allows a process to execute an executable file whose permission
    bits or ACL do not allow the process execute permission.

```

`file_dac_read`
 Allows a process to read a file or directory whose permission bits or ACL do not allow the process read permission.

`file_dac_search`
 Allows a process to search a directory whose permission bits or ACL do not allow the process search permission.

`file_dac_write`
 Allows a process to write a file or directory whose permission bits or ACL do not allow the process write permission.
 In order to write files owned by uid 0 in the absence of an effective uid of 0 ALL privileges are required.

`file_link_any`
 Allows a process to create hardlinks to files owned by a uid different from the process' effective uid.

`file_owner`
 Allows a process which is not the owner of a file to modify that file's access and modification times.
 Allows a process which is not the owner of a directory to modify that directory's access and modification times.
 Allows a process which is not the owner of a file or directory to remove or rename a file or directory whose parent directory has the ``save text image after execution'' (sticky) bit set.
 Allows a process which is not the owner of a file to mount a ``namefs'' upon that file.
 (Does not apply to setting access permission bits or ACLs.)

`file_setdac`
 Allows a process which is not the owner of a file or directory to modify that file's or directory's permission bits or ACL except for the set-uid and set-gid bits.

`file_setid`
 Allows a process to change the ownership of a file or write to a file without the set-user-ID and set-group-ID bits being cleared.
 Allows a process to set the set-group-ID bit on a file or directory whose group is not the process' effective group or one of the process' supplemental groups.
 Allows a process to set the set-user-ID bit on a file with different ownership in the presence of PRIV_FILE_SETDAC.
 Additional restrictions apply when creating or modifying a set-uid 0 file.

`ipc_dac_read`
 Allows a process to read a System V IPC Message Queue, Semaphore Set, or Shared Memory Segment whose permission bits do not allow the process read permission.
 Allows a process to read remote shared memory whose permission bits do not allow the process read permission.

`ipc_dac_write`
 Allows a process to write a System V IPC Message Queue, Semaphore Set, or Shared Memory Segment whose permission bits do not allow the process write permission.
 Allows a process to read remote shared memory whose permission bits do not allow the process write permission.
 Additional restrictions apply if the owner of the object has uid 0 and the effective uid of the current process is not 0.

`ipc_owner`

Allows a process which is not the owner of a System V IPC Message Queue, Semaphore Set, or Shared Memory Segment to remove, change ownership of, or change permission bits of the Message Queue, Semaphore Set, or Shared Memory Segment. Additional restrictions apply if the owner of the object has uid 0 and the effective uid of the current process is not 0.

`net_icmpaccess`
Allows a process to send and receive ICMP packets.

`net_privaddr`
Allows a process to bind to a privileged port number. The privilege port numbers are 1-1023 (the traditional UNIX privileged ports) as well as those ports marked as "udp/tcp_extra_priv_ports" with the exception of the ports reserved for use by NFS.

`proc_chroot`
Allows a process to change its root directory.

`proc_audit`
Allows a process to generate audit records.
Allows a process to get its own audit pre-selection information.

`proc_exec`
Allows a process to call `execve()`.

`proc_fork`
Allows a process to call `fork1()/forkall()/vfork()`

`proc_owner`
Allows a process to send signals to other processes, inspect and modify process state to other processes regardless of ownership. When modifying another process, additional restrictions apply: the effective privilege set of the attaching process must be a superset of the target process' effective, permitted and inheritable sets; the limit set must be a superset of the target's limit set; if the target process has any uid set to 0 all privilege must be asserted unless the effective uid is 0.
Allows a process to bind arbitrary processes to CPUs.

`proc_session`
Allows a process to send signals or trace processes outside its session.

`proc_setid`
Allows a process to set its uids at will.
Assuming uid 0 requires all privileges to be asserted.

`proc_taskid`
Allows a process to assign a new task ID to the calling process.

`sys_acct`
Allows a process to enable and disable and manage accounting through `acct(2)`, `getacct(2)`, `putacct(2)` and `wracct(2)`.

`sys_admin`
Allows a process to perform system administration tasks such as setting node and domain name and specifying `nscd` and `coreadm` settings.

`sys_mount`
Allows a process to mount and unmount filesystems which would otherwise be restricted (i.e., most filesystems except `namefs`)
Allows a process to add and remove swap devices.

`sys_nfs`
Allows a process to perform Sun private NFS specific system calls.

```
Allows a process to bind to ports reserved by NFS: ports 2049 (nfs)
and port 4045 (lockd).
sys_resource
Allows a process to modify the resource limits with
by setrlimit(2) and setrctl(2) without restriction.
Allows a process to exceed the per-user maximum number of
processes.
Allows a process to extend or create files on a filesystem that
has less than minfree space in reserve.
```

Mounting File Systems in Running Non-Global Zones

You can mount file systems in a running non-global zone. The following procedures are covered.

- As the global administrator in the global zone, you can import raw and block devices into a non-global zone. After the devices are imported, the zone administrator has access to the disk. The zone administrator can then create a new file system on the disk and perform one of the following actions:
 - Mount the file system manually
 - Place the file system in `/etc/vfstab` so that it will be mounted on zone boot
- As the global administrator, you can also mount a file system from the global zone into the non-global zone.

▼ How to Import Raw and Block Devices by Using `zonectfg`

This procedure uses the `lofi` command. For information about this command, see the `lofiadm(1M)` and `lofi(7D)` man pages.

1. Become superuser, or assume the Primary Administrator role.

To create the role and assign the role to a user, see “Using the Solaris Management Tools With RBAC (Task Map)” in *System Administration Guide: Basic Administration*.

2. Change directories to `/usr/tmp`.

```
global# cd /usr/tmp
```

3. Create a new UFS file system.

```
global# mkfile 10m fsfile
```

4. Attach the file as a block device.

The first available slot, which is `/dev/lofi/1` if no other `lofi` devices have been created, is used.

```
global# lofiadm -a `pwd`/fsfile
```

You will also get the required character device.

5. Import the devices into the zone `my-zone`.

```
global# zonecfg -z my-zone
zonecfg:my-zone> add device
zonecfg:my-zone:device> set match=/dev/rlofi/1
zonecfg:my-zone:device> end
zonecfg:my-zone> add device
zonecfg:my-zone:device> set match=/dev/lofi/1
zonecfg:my-zone:device> end
```

6. Reboot the zone.

```
global# zoneadm -z my-zone boot
```

7. Log in to the zone and verify that the devices were successfully imported.

```
my-zone# ls -l /dev/*lofi/*
```

You will see a display that is similar to this:

```
brw----- 1 root    sys      147, 1 Jan  7 11:26 /dev/lofi/1
crw----- 1 root    sys      147, 1 Jan  7 11:26 /dev/rlofi/1
```

▼ How to Mount the File System Manually

You must be the zone administrator and have the Zone Management profile to perform this procedure. This procedure uses the `newfs` command, which is described in the `newfs(1M)` man page.

1. Become superuser, or have the Zone Management rights profile in your list of profiles.
2. In the zone `my-zone`, create a new file system on the disk.

```
my-zone# newfs /dev/lofi/1
```

3. Respond yes at the prompt.

```
newfs: construct a new file system /dev/rlofi/1: (y/n)? y
```

You will see a display that is similar to this:

```
/dev/rlofi/1: 20468 sectors in 34 cylinders of 1 tracks, 602 sectors
              10.0MB in 3 cyl groups (16 c/g, 4.70MB/g, 2240 i/g)
super-block backups (for fsck -F ufs -o b=#) at:
              32, 9664, 19296,
```


4. Check the file system for errors.

```
my-zone# fsck -F ufs /dev/rlofi/1
```

You will see a display that is similar to this:

```
** /dev/rlofi/1
** Last Mounted on
** Phase 1 - Check Blocks and Sizes
** Phase 2 - Check Pathnames
** Phase 3 - Check Connectivity
** Phase 4 - Check Reference Counts
** Phase 5 - Check Cyl groups
2 files, 9 used, 9320 free (16 frags, 1163 blocks, 0.2% fragmentation)
```

5. Mount the file system.

```
my-zone# mount -F ufs /dev/lofi/1 /mnt
```

6. Verify the mount.

```
my-zone# grep /mnt /etc/mnttab
```

You will see a display similar to this:

```
/dev/lofi/1 /mnt ufs
rw,suid,intr,largefiles,xattr,onerror=panic,zone=foo,dev=24c0001
1073503869
```

▼ How to Place a File System in `/etc/vfstab` to Be Mounted When the Zone Boots

This procedure is used to mount the block device `/dev/lofi/1` on the file system path `/mnt`. The block device contains a UFS file system. The following options are used:

- `logging` is used as the mount option.
- `yes` tells the system to automatically mount the file system when the zone boots.
- `/dev/rlofi/1` is the character (or raw) device. The `fsck` command is run on the raw device if required.

1. Become superuser, or have the Zone Management rights profile in your list of profiles.

2. In the zone `my-zone`, add the following line to `/etc/vfstab`:

```
/dev/lofi/1 /dev/rlofi/1 /mnt ufs 2 yes logging
```

▼ How to Mount a File System From the Global Zone Into a Non-Global Zone

Assume that a zone has the `zonepath /export/home/my-zone`. You want to mount the disk `/dev/lofi/1` from the global zone into `/mnt` in the non-global zone.

For information about `lofi`, see the `lofiadm(1M)` and `lofi(7D)` man pages.

You must be the global administrator in the global zone to perform this procedure.

1. Become superuser, or assume the Primary Administrator role.

To create the role and assign the role to a user, see “Using the Solaris Management Tools With RBAC (Task Map)” in *System Administration Guide: Basic Administration*.

2. To mount the disk into `/mnt` in the non-global zone, type the following from the global zone:

```
global# mount -F ufs /dev/lofi/1 /export/home/my-zone/root/mnt
```

Using IP Network Multipathing in a Zones Environment

▼ How to Extend IP Network Multipathing Functionality to Non-Global Zones

Use this procedure to configure IP Network Multipathing (IPMP) in the global zone and extend its functionality to non-global zones.

Each address, or logical interface, should be associated with a non-global zone when you configure the zone. See “Using the `zonecfg` Command” on page 209 and “How to Configure the Zone” on page 226 for instructions.

This procedure accomplishes the following:

- The cards `bge0` and `hme0` are configured together in a group.
- Address `192.168.0.1` is associated with the non-global zone `my-zone`.
- The `bge0` card is set as the physical interface. Thus, the IP address is hosted in the group that contains the `bge0` and `hme0` cards.

In a running zone, you can use the `ifconfig` command to make the association. See “Network Interfaces” on page 301 and the `ifconfig(1M)` man page.

You must be the global administrator in the global zone to perform this procedure.

1. Become superuser, or assume the Primary Administrator role.

To create the role and assign the role to a user, see “Using the Solaris Management Tools With RBAC (Task Map)” in *System Administration Guide: Basic Administration*.

2. In the global zone, configure IPMP groups as described in “Configuring Multipathing Interface Groups” in *System Administration Guide: IP Services*.

3. Use the `zonectfg` command to configure the zone. When you configure the `net` resource, add address `192.168.0.1` and physical interface `bge0` to the zone `my-zone`:

```
zonectfg:my-zone> add net
zonectfg:my-zone:net> set address=192.168.0.1
zonectfg:my-zone:net> set physical=bge0
zonectfg:my-zone:net> end
```

Only `bge0` would be visible in non-global zone `my-zone`.

If `bge0` subsequently fails and the `bge0` data addresses fail over to `hme0` in the global zone, then the `my-zone` addresses migrate as well.

If address `192.168.0.1` moves to `hme0`, then only `hme0` would now be visible in non-global zone `my-zone`. This card would be associated with address `192.168.0.1`, and `bge0` would no longer be visible.

Using the Fair Share Scheduler in a Zones Environment

How to Set FSS Shares in the Global Zone

The global zone is given one share by default. You can use this procedure to change the default allocation.

You must be the global administrator in the global zone to perform this procedure.

1. Become superuser, or assume the Primary Administrator role.

To create the role and assign the role to a user, see “Using the Solaris Management Tools With RBAC (Task Map)” in *System Administration Guide: Basic Administration*.

2. Use the `prctl` utility to assign two shares to the global zone:

```
# prctl -n zone.cpu-shares -v 2 -r -i zone global
```

Note that you must reset shares allocated through the `prctl` utility whenever you reboot the system.

3. (Optional) To verify the number of shares assigned to the global zone, type:

```
# prctl -n zone.cpu-shares -i zone global
```

For more information on the `prctl` utility, see the `prctl(1)` man page.

▼ How to Balance CPU Usage Between the Global Zone and Non-Global Zones

This procedure can be used to equalize the CPU usage between the global zone and non-global zones. In this procedure, a configuration with two non-global zones is used. The global zone receives one-third of the machine's CPU resource. The non-global zones each receive one-third of the machine's CPU resource.

You must be the global administrator in the global zone to perform this procedure.

1. **Become superuser, or assume the Primary Administrator role.**

To create the role and assign the role to a user, see "Using the Solaris Management Tools With RBAC (Task Map)" in *System Administration Guide: Basic Administration*.

2. **Use the `zonecfg` command to assign one share to each non-global zone as shown in [Step 11](#).**

Because the global zone receives one share by default, also giving each non-global zone one share divides the machine equally.

Using Rights Profiles in Zone Administration

This section covers tasks associated with using rights profiles in non-global zones.

▼ How to Assign the Zone Management Profile

The Zone Management profile grants the power to manage all of the non-global zones on the system to a user.

You must be the global administrator in the global zone to perform this procedure.

1. **Become superuser, or assume the Primary Administrator role.**

To create the role and assign the role to a user, see “Using the Solaris Management Tools With RBAC (Task Map)” in *System Administration Guide: Basic Administration*.

2. **Create a role that includes the Zone Management rights profile, and assign the role to a user.**

- To create and assign the role by using the Solaris Management Console, see “Configuring RBAC (Task Map)” in *System Administration Guide: Security Services*. Refer to the task “How to Create and Assign a Role By Using the GUI.”
- To create and assign the role on the command line, see “Managing RBAC” in *System Administration Guide: Security Services*. Refer to the task “How to Create a Role From the Command Line.”

Example—Using Profile Shells With Zone Commands

You can execute zone commands in a profile using the `pfexec` program. The program executes commands with the attributes specified by the user’s profiles in the `exec_attr` database. The program is invoked by the profile shells `pfksh`, `pfcsch`, and `pfsh`.

Use the `pfexec` program to log in to a zone, for example, `my-zone`.

```
machine$ pfexec zlogin my-zone
```


Glossary

bless	In Perl, the keyword used to create an object.
blessed	In Perl, the term used to denote class membership. <i>bless</i>
cap	A limit that is placed on system resource usage.
capping	The process of placing a limit on system resource usage.
default pool	The pool created by the system when pools are enabled. See also resource pool .
default processor set	The processor set created by the system when pools are enabled. See also processor set .
disjoint	A type of set in which the members of the set do not overlap and are not duplicated.
dynamic configuration	Information about the disposition of resources within the resource pools framework for a given system at a point in time.
dynamic reconfiguration	On SPARC based systems, the ability to reconfigure hardware while the system is running. Also known as DR.
extended accounting	A flexible way to record resource consumption on a task basis or process basis in the Solaris Operating System.
fair share scheduler	A scheduling class, also known as FSS, that allows you to allocate CPU time that is based on shares. Shares define the portion of the system's CPU resources allocated to a project.
FSS	See fair share scheduler .
global administrator	An administrator with superuser privileges or the Primary Administrator role. When logged in to the global zone, the global administrator can monitor and control the system as a whole. See also zone administrator .

global zone	The zone contained on every Solaris system. In a zones environment, the global zone is both the default zone for the system and the zone used for system-wide administrative control. See also non-global zone .
heap	Process-allocated scratch memory.
locked memory	Memory that cannot be paged.
memory cap enforcement threshold	The percentage of physical memory utilization on the system that will trigger cap enforcement by the resource capping daemon.
naming service database	In the “Projects and Tasks (Overview)” chapter of this document, a reference to both LDAP containers and NIS maps.
non-global zone	A virtualized operating system environment created within a single instance of the Solaris Operating System. The Solaris Zones software partitioning technology is used to virtualize operating system services.
non-global zone administrator	See zone administrator .
page in	To read data from a file into physical memory one page at a time.
page out	To relocate pages to an area outside of physical memory.
pool	See resource pool .
pool daemon	The <code>poold</code> system daemon that is active when dynamic resource allocation is required.
project	A network-wide administrative identifier for related work.
processor set	A disjoint grouping of CPUs. Each processor set can contain zero or more processors. A processor set is represented in the resource pools configuration as a resource element. Also referred to as a <code>pset</code> . See also disjoint .
resident set size	The size of the resident set. The resident set is the set of pages that are resident in physical memory.
resource	An aspect of the computing system that can be manipulated with the intent to change application behavior.
resource capping daemon	A daemon that regulates the consumption of physical memory by processes running in projects that have resource caps defined.
resource consumer	Fundamentally, a Solaris process. Process model entities such as the project and the task provide ways of discussing resource consumption in terms of aggregated resource consumption.
resource control	A per-process, per-task, or per-project limit on the consumption of a resource.

resource management	A functionality that enables you to control how applications use available system resources.
resource partition	An exclusive subset of a resource. All of the partitions of a resource sum to represent the total amount of the resource available in a single executing Solaris instance.
resource pool	A configuration mechanism that is used to partition machine resources. A resource pool represents an association between groups of partitionable resources.
resource set	A process-bindable resource. Most often used to refer to the objects constructed by a kernel subsystem offering some form of partitioning. Examples of resource sets include scheduling classes and processor sets.
RSS	See resident set size .
scanner	A kernel thread that identifies infrequently used pages and relocates the pages to an area outside of physical memory.
Solaris Zones	A software partitioning technology used to virtualize operating system services and provide an isolated, secure environment in which to run applications.
static pools configuration	A representation of the way in which an administrator would like a system to be configured with respect to resource pools functionality.
task	In resource management, a process collective that represents a set of work over time. Each task is associated with one project.
working set size	The size of the working set. The working set is the set of pages that the project workload actively uses during its processing cycle.
workload	An aggregation of all processes of an application or group of applications.
WSS	See working set size .
zone administrator	An administrator having the Zone Management profile. The privileges of a zone administrator are confined to a non-global zone. See also global administrator .
zone state	The status of a non-global zone. The zone state is one of configured, incomplete, installed, ready, running, or shutting down.

Index

A

acctadm command, 65, 66
activating extended accounting, 64-66
administering resource pools, 155
attribute, `project.pool`, 135

B

binding to resource pool, 171

C

changing resource controls temporarily, 83
commands
 extended accounting, 59
 fair share scheduler, 106
 projects and tasks, 40
 resource controls, 84
 zones, 310
configuration, `rcapd`, 115
configuring resource controls, 73
configuring zones, tasks, 219
CPU share configuration, 100
creating resource pools, 140

D

deactivating extended accounting, 66
default processor set, 134
default project, 34

default resource pool, 134
disabling resource capping, 126
displaying extended accounting status, 65
DRP, 134
dynamic pools configuration, 137

E

enabling resource capping, 125
entry format, `/etc/project` file, 36
 `/etc/project`
 entry format, 36
 file, 35
 `/etc/user_attr` file, 34
exact file, 56
extended accounting
 activating, 64-66
 chargeback, 56
 commands, 59
 deactivating, 66
 file format, 56
 overview, 55
 status, displaying, 65

F

fair share scheduler
 and processor sets, 102
 `project.cpu-shares`, 96
 share definition, 96
fair share scheduler (FSS), 96

FSS

See fair share scheduler (FSS)
configuration, 109

G

global administrator, 198, 200
global zone, 198

H

halting a zone, 241
troubleshooting, 241

I

implementing resource pools, 139
installing zones, 244
interactive packages, 267
interprocess communication, *See* resource
controls
IPC, 72

L

libexecct, 56
login, remote zone, 256

M

memory cap enforcement threshold, 116

N

non-global zone, 198

P

package operations, 267
packages, interactive, 267

PAM (pluggable authentication module),
identity management, 36
Perl interface, 59
pluggable authentication module, *See* PAM
pool, asynchronous control violation, 152
pool
configurable features, 147
constraints, 143
control scope, 152
cpu-pinned property, 143
description, 141
dynamic resource allocation, 134
logging information, 148
objectives, 143
synchronous control violation, 152
pools, 133
poolstat
description, 153
output format, 154
usage examples, 173
populating a zone, 238
privilege levels, 78
project
active state, 97
definition, 34
idle state, 97
with zero shares, 96
project 0, 101
project.cpu-shares, 100
project database, 35
project.pool attribute, 135
project system, *See* project 0
putacct, 57

R

rcap.max-rss, 115
rcapadm, 115
rcapd, 113
sample intervals, 119
scan intervals, 118
rcapd configuration, 115
rcapstat, 119
rctls, 71
See resource controls
remote zone login, 256
removing resource pools, 171

- resource cap, 113
- resource capping
 - disabling, 126
 - enabling, 125
- resource capping daemon, 113
- resource controls
 - changing temporarily, 83
 - configuring, 73
 - definition, 71
 - global actions, 79
 - inf value, 82
 - interprocess communication, 72
 - list of, 74
 - local actions, 79
 - overview, 71
 - temporarily updating, 83
 - threshold values, 79
 - zone-wide, 216
- resource limits, 72
- resource management
 - constraints, 28
 - definition, 26
 - partitioning, 28
 - scheduling, 28
- resource pools, 133
 - activating configuration, 170
 - administering, 155
 - binding to, 171
 - configuration elements, 138
 - creating, 140
 - disabling, 159
 - dynamic reconfiguration, 140
 - enabling, 159
 - `/etc/pooladm.conf`, 138
 - implementing, 139
 - properties, 138
 - removing, 171
 - removing configuration, 170
 - static pools configuration, 137
- rlimits, *See* resource limits

S

- server consolidation, 29
- Solaris Management Console
 - definition, 184
 - performance monitoring, 185

- Solaris Management Console (Continued)
 - setting resource controls, 191
 - SUNW_PKG_ALLZONES package parameter, 274
 - SUNW_PKG_HOLLOW package parameter, 276

T

- tasks, resource management, 39
- temporarily updating resource controls, 83
- threshold values, 78

V

- `/var/adm/exacct` directory, 58

Z

- zone
 - adding packages, 269
 - adding patches, 279
 - boot procedure, 246
 - commands used in, 310
 - configuration, 209
 - creating, 200
 - definition, 195
 - disk space, 221
 - features, 203
 - halt procedure, 249
 - halting, 241
 - installation, 245
 - interactive mode, 256
 - network address, 223
 - non-interactive mode, 256
 - package and patch overview, 265
 - package rules, 267
 - PatchPro support, 281
 - populating, 238
 - reboot, 241
 - reboot procedure, 250
 - removing packages, 272
 - removing patches, 280
 - resource controls, 216
 - resource type properties, 214
 - resource types, 213
 - scope, 265

- zone (Continued)
 - states, 201
 - uninstall procedure, 250
- zone administrator, 200
- zone configuration
 - overview, 207
 - script, 230
 - tasks, 219
- zone console login, console login mode, 255
- zone.cpu-shares, 216
 - zone resource control, 212
- zone host name, 223
- zone ID, 198
- zone installation
 - overview, 237
 - tasks, 244
- zone login
 - failsafe mode, 255
 - overview, 253
 - remote, 256
- zone.max-lwps, 216
 - zone resource control, 212
- zone name, 198
- zone node name, 295
- zone resource controls, 212
- zone size, restricting, 222
- zone-wide resource controls, 216
- zoneadmd, 239
- zonecfg
 - entities, 212
 - modes, 210
 - operations, 207
 - procedure, 226
 - scope, 210
 - scope, global, 210
 - scope, resource specific, 210
 - subcommands, 210
- zones, characteristics by type, 199
- zones commands, 310
- zsched, 240