

# Oracle® Solaris Tunable Parameters Reference Manual

Copyright © 2000, 2010, Oracle and/or its affiliates. All rights reserved.

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

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

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

U.S. GOVERNMENT RIGHTS Programs, software, databases, and related documentation and technical data delivered to U.S. Government customers are "commercial computer software" or "commercial technical data" pursuant to the applicable Federal Acquisition Regulation and agency-specific supplemental regulations. As such, the use, duplication, disclosure, modification, and adaptation shall be subject to the restrictions and license terms set forth in the applicable Government contract, and, to the extent applicable by the terms of the Government contract, the additional rights set forth in FAR 52.227-19, Commercial Computer Software License (December 2007). Oracle America, Inc., 500 Oracle Parkway, Redwood City, CA 94065.

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

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

AMD, Opteron, the AMD logo, and the AMD Opteron logo are trademarks or registered trademarks of Advanced Micro Devices. Intel and Intel Xeon are trademarks or registered trademarks of Intel Corporation. All SPARC trademarks are used under license and are trademarks or registered trademarks of SPARC International, Inc. UNIX is a registered trademark licensed through X/Open Company, Ltd.

This software or hardware and documentation may provide access to or information on content, products, and services from third parties. Oracle Corporation and its affiliates are not responsible for and expressly disclaim all warranties of any kind with respect to third-party content, products, and services. Oracle Corporation and its affiliates will not be responsible for any loss, costs, or damages incurred due to your access to or use of third-party content, products, or services.

---

Copyright © 2000, 2010, Oracle et/ou ses affiliés. Tous droits réservés.

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

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

Si ce logiciel, ou la documentation qui l'accompagne, est concédé sous licence au Gouvernement des Etats-Unis, ou à toute entité qui délivre la licence de ce logiciel ou l'utilise pour le compte du Gouvernement des Etats-Unis, la notice suivante s'applique :

U.S. GOVERNMENT RIGHTS Programs, software, databases, and related documentation and technical data delivered to U.S. Government customers are "commercial computer software" or "commercial technical data" pursuant to the applicable Federal Acquisition Regulation and agency-specific supplemental regulations. As such, the use, duplication, disclosure, modification, and adaptation shall be subject to the restrictions and license terms set forth in the applicable Government contract, and, to the extent applicable by the terms of the Government contract, the additional rights set forth in FAR 52.227-19, Commercial Computer Software License (December 2007). Oracle America, Inc., 500 Oracle Parkway, Redwood City, CA 94065.

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

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

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

# Contents

---

<b>Preface</b> .....	13
<b>1 Overview of Oracle Solaris System Tuning</b> .....	17
What's New in Oracle Solaris System Tuning? .....	17
Tuning a Solaris System .....	19
Tuning Format of Tunable Parameters Descriptions .....	20
Tuning the Solaris Kernel .....	21
/etc/system File .....	22
kldb Command .....	23
mdb Command .....	23
Special Solaris tune and var Structures .....	24
Viewing Solaris System Configuration Information .....	24
sysdef Command .....	24
kstat Utility .....	25
<b>2 Oracle Solaris Kernel Tunable Parameters</b> .....	27
Where to Find Tunable Parameter Information .....	28
General Kernel and Memory Parameters .....	28
physmem .....	28
zfs_arc_min .....	29
zfs_arc_max .....	29
default_stksize .....	30
lwp_default_stksize .....	31
logevent_max_q_sz .....	32
segkpsize .....	33
noexec_user_stack .....	33
fsflush and Related Parameters .....	34

fsflush .....	34
tune_t_fsflushr .....	35
autoup .....	36
dopageflush .....	37
doiflush .....	37
Process-Sizing Parameters .....	38
maxusers .....	38
reserved_procs .....	39
pidmax .....	40
max_nprocs .....	41
maxuprc .....	41
ngroups_max .....	42
Paging-Related Parameters .....	42
lotsfree .....	44
desfree .....	45
minfree .....	46
throttlefree .....	47
pageout_reserve .....	47
pages_pp_maximum .....	48
tune_t_minarmem .....	49
fastscan .....	50
slowscan .....	50
min_percent_cpu .....	51
handspreadpages .....	51
pages_before_pager .....	52
maxpgio .....	53
Swapping-Related Parameters .....	54
swapfs_reserve .....	54
swapfs_minfree .....	54
Kernel Memory Allocator .....	55
kmem_flags .....	55
kmem_stackinfo .....	57
General Driver Parameters .....	58
moddebug .....	58
ddi_msix_alloc_limit .....	59
Network Driver Parameters .....	60

---

igb Parameters .....	60
ixgbe Parameters .....	61
General I/O Parameters .....	65
maxphys .....	65
rlim_fd_max .....	65
rlim_fd_cur .....	66
General File System Parameters .....	67
ncsize .....	67
rstchown .....	68
dnlc_dir_enable .....	69
dnlc_dir_min_size .....	69
dnlc_dir_max_size .....	70
segmap_percent .....	70
UFS Parameters .....	71
bufhwm and bufhwm_pct .....	71
ndquot .....	73
ufs_ninode .....	73
ufs_WRITES .....	75
ufs_LWand ufs_HW .....	75
freebehind .....	76
smallfile .....	77
TMPFS Parameters .....	78
tmpfs:tmpfs_maxkmem .....	78
tmpfs:tmpfs_minfree .....	78
Pseudo Terminals .....	79
pt_cnt .....	80
pt_pctofmem .....	80
pt_max_pty .....	81
STREAMS Parameters .....	82
nstrpush .....	82
strmsgsz .....	82
strctlsz .....	83
System V Message Queues .....	83
System V Semaphores .....	83
System V Shared Memory .....	84
segspt_minfree .....	84

Scheduling .....	85
rechoose_interval .....	85
Timers .....	85
hires_tick .....	85
timer_max .....	86
sun4u or sun4v Specific Parameters .....	86
consistent_coloring .....	86
tsb_alloc_hiwater_factor .....	87
default_tsb_size .....	88
enable_tsb_rss_sizing .....	89
tsb_rss_factor .....	89
Locality Group Parameters .....	90
lpg_alloc_prefer .....	90
lgrp_mem_default_policy .....	91
lgrp_mem_pset_aware .....	92
<b>3 NFS Tunable Parameters .....</b>	<b>95</b>
Where to Find Tunable Parameter Information .....	95
Tuning the NFS Environment .....	95
NFS Module Parameters .....	96
nfs:nfs3_pathconf_disable_cache .....	96
nfs:nfs4_pathconf_disable_cache .....	96
nfs:nfs_allow_preepoch_time .....	97
nfs:nfs_cots_timeo .....	98
nfs:nfs3_cots_timeo .....	98
nfs:nfs4_cots_timeo .....	99
nfs:nfs_do_symlink_cache .....	100
nfs:nfs3_do_symlink_cache .....	100
nfs:nfs4_do_symlink_cache .....	101
nfs:nfs_dynamic .....	102
nfs:nfs3_dynamic .....	102
nfs:nfs_lookup_neg_cache .....	103
nfs:nfs3_lookup_neg_cache .....	103
nfs:nfs4_lookup_neg_cache .....	104
nfs:nfs_max_threads .....	105

---

nfs:nfs3_max_threads .....	106
nfs:nfs4_max_threads .....	107
nfs:nfs_nra .....	107
nfs:nfs3_nra .....	108
nfs:nfs4_nra .....	109
nfs:nrnode .....	109
nfs:nfs_shrinkreaddir .....	110
nfs:nfs3_shrinkreaddir .....	111
nfs:nfs_write_error_interval .....	112
nfs:nfs_write_error_to_cons_only .....	112
nfs:nfs_disable_rmdir_cache .....	113
nfs:nfs_bsize .....	114
nfs:nfs3_bsize .....	114
nfs:nfs4_bsize .....	115
nfs:nfs_async_clusters .....	116
nfs:nfs3_async_clusters .....	116
nfs:nfs4_async_clusters .....	117
nfs:nfs_async_timeout .....	118
nfs:nacache .....	119
nfs:nfs3_jukebox_delay .....	120
nfs:nfs3_max_transfer_size .....	120
nfs:nfs4_max_transfer_size .....	121
nfs:nfs3_max_transfer_size_clts .....	122
nfs:nfs3_max_transfer_size_cots .....	123
nfssrv Module Parameters .....	123
nfssrv:nfs_portmon .....	123
nfssrv:rfs_write_async .....	124
nfssrv:nfsauth_ch_cache_max .....	125
nfssrv:exi_cache_time .....	126
rpcmod Module Parameters .....	126
rpcmod:clnt_max_conns .....	126
rpcmod:clnt_idle_timeout .....	127
rpcmod:svc_idle_timeout .....	127
rpcmod:svc_default_stksize .....	128
rpcmod:svc_default_max_same_xprt .....	129
rpcmod:maxdupreqs .....	129

rpcmod:cotsmaxdupreqs .....	130
<b>4 Internet Protocol Suite Tunable Parameters .....</b>	<b>133</b>
Where to Find Tunable Parameter Information .....	133
Overview of Tuning IP Suite Parameters .....	133
IP Suite Parameter Validation .....	134
Internet Request for Comments (RFCs) .....	134
IP Tunable Parameters .....	134
_icmp_err_interval and_icmp_err_burst .....	134
_respond_to_echo_broadcast .....	135
_addrs_per_if .....	135
ip_queue_fanout .....	136
TCP Tunable Parameters .....	136
_deferred_ack_interval .....	136
_local_dack_interval .....	137
_deferred_acks_max .....	137
_local_dacks_max .....	138
_wscale_always .....	138
_tstamp_always .....	139
send_maxbuf .....	139
recv_maxbuf .....	139
_max_buf .....	140
_cwnd_max .....	140
_slow_start_initial .....	141
_slow_start_after_idle .....	141
sack .....	141
_rev_src_routes .....	142
_time_wait_interval .....	142
ecn .....	143
_conn_req_max_q .....	144
_conn_req_max_q0 .....	144
_conn_req_min .....	145
_rst_sent_rate_enabled .....	145
_rst_sent_rate .....	146
TCP/IP Parameters Set in the /etc/system File .....	146



TCP Parameters With Additional Cautions .....	147
UDP Tunable Parameters .....	151
send_maxbuf .....	151
recv_maxbuf .....	151
IPQoS Tunable Parameter .....	152
_policy_mask .....	152
SCTP Tunable Parameters .....	153
_max_init_retr .....	153
_pa_max_retr .....	153
_pp_max_retr .....	154
_cwnd_max .....	154
_ipv4_ttl .....	154
_heartbeat_interval .....	155
_new_secret_interval .....	155
_initial_mtu .....	155
_deferred_ack_interval .....	156
_ignore_path_mtu .....	156
_initial_ssthresh .....	156
_max_buf .....	157
_ipv6_hoplimit .....	157
_rto_min .....	157
_rto_max .....	158
_rto_initial .....	158
_cookie_life .....	158
_max_in_streams .....	158
_initial_out_streams .....	159
_shutack_wait_bound .....	159
_maxburst .....	159
_addip_enabled .....	160
_prsctp_enabled .....	160
Per-Route Metrics .....	160
<b>5 Network Cache and Accelerator Tunable Parameters .....</b>	<b>163</b>
Where to Find Tunable Parameters Information .....	163
Tuning NCA Parameters .....	163

nca:nca_conn_hash_size .....	164
nca:nca_conn_req_max_q .....	164
nca:nca_conn_req_max_q0 .....	164
nca:nca_ppmax .....	165
nca:nca_vpmax .....	165
General System Tuning for the NCA .....	166
sq_max_size .....	166
ge:ge_intr_mode .....	167
<b>6 System Facility Parameters .....</b>	<b>169</b>
System Default Parameters .....	170
autofs .....	170
cron .....	170
devfsadm .....	170
dhcpgent .....	170
fs .....	170
ftp .....	171
inetinit .....	171
init .....	171
ipsec .....	171
kbd .....	171
keyserv .....	171
login .....	171
mpathd .....	172
nfs .....	172
nfslogd .....	172
nss .....	172
passwd .....	172
power .....	172
su .....	172
syslog .....	172
sys-suspend .....	173
tar .....	173
utmpd .....	173
yppasswdd .....	173

---

<b>A Tunable Parameters Change History</b> .....	175
Kernel Parameters .....	175
Process-Sizing Tunables .....	175
General Driver Parameter .....	175
Network Driver Parameters .....	176
General Kernel and Memory Parameters .....	176
fsflush and Related Parameters .....	176
Parameters That Are Obsolete or Have Been Removed .....	176
TCP/IP Module Parameters .....	176
<b>B Revision History for This Manual</b> .....	179
Current Version: Oracle Solaris 11 Express Release .....	179
New or Changed Parameters in the Oracle Solaris Release .....	179
<b>Index</b> .....	181



# Preface

---

The *Oracle Solaris Tunable Parameters Reference Manual* provides reference information about Oracle Solaris OS kernel and network tunable parameters. This manual does not provide tunable parameter information about desktop systems or Java environments.

This manual contains information for both SPARC based and x86 based systems.

---

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

In this document these x86 terms mean the following:

- “x86” refers to the larger family of 64-bit and 32-bit x86 compatible products.
- “x64” relates specifically to 64-bit x86 compatible CPUs.
- “32-bit x86” points out specific 32-bit information about x86 based systems.

For supported systems, see *Oracle Solaris Hardware Compatibility List* at <http://www.sun.com/bigadmin/hcl>.

---

## Who Should Use This Book

This book is intended for experienced Solaris system administrators who might need to change kernel tunable parameters in certain situations. For guidelines on changing Solaris tunable parameters, refer to “[Tuning a Solaris System](#)” on page 19.

## How This Book Is Organized

The following table describes the chapters and appendixes in this book.

---

Chapter	Description
<a href="#">Chapter 1, “Overview of Oracle Solaris System Tuning”</a>	An overview of tuning a Solaris system. Also provides a description of the format used in the book to describe the kernel tunables.

---

Chapter	Description
Chapter 2, “Oracle Solaris Kernel Tunable Parameters”	A description of Solaris kernel tunables such as kernel memory, file system, process size, and paging parameters.
Chapter 3, “NFS Tunable Parameters”	A description of NFS tunables such as caching symbolic links, dynamic retransmission, and RPC security parameters.
Chapter 4, “Internet Protocol Suite Tunable Parameters”	A description of TCP/IP tunables such as IP forwarding, source routing, and buffer-sizing parameters.
Chapter 5, “Network Cache and Accelerator Tunable Parameters”	A description of tunable parameters for the Network Cache and Accelerator (NCA).
Chapter 6, “System Facility Parameters”	A description of parameters used to set default values of certain system facilities. Changes are made by modifying files in the <code>/etc/default</code> directory.
Appendix A, “Tunable Parameters Change History”	A history of parameters that have changed or are now obsolete.
Appendix B, “Revision History for This Manual”	A history of this manual's revisions including the current Solaris release.

## Other Resources for Solaris Tuning Information

This table describes other resources for Solaris tuning information.

Tuning Resource	For More Information
Online performance tuning information	<a href="http://www.solarisinternals.com/si/index.php">http://www.solarisinternals.com/si/index.php</a>
In-depth technical white papers	<a href="http://developers.sun.com/solaris/">http://developers.sun.com/solaris/</a>

## Documentation, Support, and Training

See the following web sites for additional resources:

- Documentation (<http://www.oracle.com/technetwork/indexes/documentation/index.html>)
- Support (<http://www.oracle.com/us/support/systems/index.html>)
- Training (<http://education.oracle.com>) – Click the Sun link in the left navigation bar.

## Oracle Software Resources

Oracle Technology Network (<http://www.oracle.com/technetwork/index.html>) offers a range of resources related to Oracle software:

- Discuss technical problems and solutions on the [Discussion Forums](http://forums.oracle.com) (<http://forums.oracle.com>).
- Get hands-on step-by-step tutorials with [Oracle By Example](http://www.oracle.com/technetwork/tutorials/index.html) (<http://www.oracle.com/technetwork/tutorials/index.html>).
- Download [Sample Code](http://www.oracle.com/technology/sample_code/index.html) ([http://www.oracle.com/technology/sample\\_code/index.html](http://www.oracle.com/technology/sample_code/index.html)).

## Typographic Conventions

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

TABLE P-1 Typographic Conventions

Typeface	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% you have mail.</code>
<b>AaBbCc123</b>	What you type, contrasted with onscreen computer output	<code>machine_name% su</code> Password:
<i>aabbcc123</i>	Placeholder: replace with a real name or value	The command to remove a file is <i>rm filename</i> .
<i>AaBbCc123</i>	Book titles, new terms, and terms to be emphasized	Read <i>Chapter 6</i> in the <i>User's Guide</i> . <i>A cache</i> is a copy that is stored locally. Do <i>not</i> save the file. <b>Note:</b> Some emphasized items appear bold online.

## Shell Prompts in Command Examples

The following table shows the default UNIX system prompt and superuser prompt for shells that are included in the Oracle Solaris OS. Note that the default system prompt that is displayed in command examples varies, depending on the Oracle Solaris release.

TABLE P-2 Shell Prompts

Shell	Prompt
Bash shell, Korn shell, and Bourne shell	\$
Bash shell, Korn shell, and Bourne shell for superuser	#
C shell	machine_name%
C shell for superuser	machine_name#



# Overview of Oracle Solaris System Tuning

---

This section provides overview information about the format of the tuning information in this manual. This section also describes the different ways to tune a Solaris system.

- “What's New in Oracle Solaris System Tuning?” on page 17
- “Tuning a Solaris System” on page 19
- “Tuning Format of Tunable Parameters Descriptions” on page 20
- “Tuning the Solaris Kernel” on page 21
- “Special Solaris tune and var Structures” on page 24
- “Viewing Solaris System Configuration Information” on page 24
- “kstat Utility” on page 25

## What's New in Oracle Solaris System Tuning?

This section describes new or changed parameters in the Oracle Solaris 11 Express release.

- The `ipadm` command replaces the `nnd` command for setting TCP/IP properties. TCP, IP, UDP, and SCTP properties are set as follows:
  - Display or set an IP property:

```
# ipadm set-prop -p property-name ipv4
# ipadm set-prop -p property-name ipv6
# ipadm show-prop -p property-name ipv4
# ipadm show-prop -p property-name ipv6
```
  - Display or set a TCP property:

```
# ipadm set-prop -p property-name tcp
# ipadm show-prop -p property-name tcp
```
  - Display or set a UDP property:

```
# ipadm set-prop -p property-name udp
# ipadm show-prop -p property-name udp
```
  - Display or set a SCTP property:

```
# ipadm set-prop -p property-name sctp
# ipadm show-prop -p property-name sctp
```

For more information, see “Overview of Tuning IP Suite Parameters” on page 133.

- This release includes the `ngroups_max` parameter description. For more information, see “`ngroups_max`” on page 42.
- This release includes the `zfs_arc_min` and `zfs_arc_max` parameter descriptions. For more information, see “`zfs_arc_min`” on page 29 and “`zfs_arc_max`” on page 29.
- This release includes several `igb` and `ixgbe` network driver parameters. For more information, see “`igb` Parameters” on page 60 and “`ixgbe` Parameters” on page 61.
- This release includes the `ddi_msix_alloc_limit` parameter that can be used to increase the number of MSI-X interrupts that a device instance can allocate. For more information, see “`ddi_msix_alloc_limit`” on page 59.
- A previous version of this manual incorrectly identified the range of the `tcp_local_dack_interval` parameter as 1 millisecond to 1 minute. The correct range is 10 milliseconds to 1 minute. For more information, see “`_local_dack_interval`” on page 137.
- This release includes the `kmem_stackinfo` parameter, which can be enabled to monitor kernel thread stack usage. For more information, see “`kmem_stackinfo`” on page 57.
- For information about tuning ZFS file systems, see the following site:  
[http://www.solarisinternals.com/wiki/index.php/ZFS\\_Evil\\_Tuning\\_Guide](http://www.solarisinternals.com/wiki/index.php/ZFS_Evil_Tuning_Guide)
- Memory locality group parameters are provided in this release. For more information about these parameters, see “`Locality Group Parameters`” on page 90.
- Parameter information was updated to include sun4v systems. For more information, see the following references:
  - “`maxphys`” on page 65
  - “`tmpfs:tmpfs_maxkmem`” on page 78
  - “`sun4u or sun4v Specific Parameters`” on page 86
- The IP instances project enables you to configure a zone as an exclusive-IP zone and assign exclusive access of some LANs or VLANs to that zone.

The previous behavior of shared-IP zones remains the default behavior. The exclusive-IP zone means that all aspects of the TCP/IP state and policy are per exclusive-IP zone, including TCP/IP tunable parameters.

The introduction of the IP instances feature means that the following TCP parameters can only be set in the global zone because they require the `PRIV_SYS_NET_CONFIG` privilege:

- “`ip_queue_fanout`” on page 136
- “`ip_queue_worker_wait`” on page 147

The other TCP, IP, and SCTP parameters and route metrics only require the `PRIV_SYS_IP_CONFIG` privilege. Each exclusive-IP zone controls its own set of these

parameters. For shared-IP zones, TCP, IP, SCTP, and route parameters are controlled by the global zone since the settings of these parameters are shared between the global zone and all shared IP zones.

For more information about using IP instances in Solaris zones, see *System Administration Guide: Oracle Solaris Zones, Oracle Solaris 10 Containers, and Resource Management*.

## Tuning a Solaris System

The Solaris OS is a multi-threaded, scalable UNIX operating system that runs on SPARC and x86 processors. It is self-adjusting to system load and demands minimal tuning. In some cases, however, tuning is necessary. This book provides details about the officially supported kernel tuning options available for the Solaris OS.

The Solaris kernel is composed of a core portion, which is always loaded, and a number of loadable modules that are loaded as references are made to them. Many variables referred to in the kernel portion of this guide are in the core portion. However, a few variables are located in loadable modules.

A key consideration in system tuning is that setting system parameters (or system variables) is often the least effective action that can be done to improve performance. Changing the behavior of the application is generally the most effective tuning aid available. Adding more physical memory and balancing disk I/O patterns are also useful. In a few rare cases, changing one of the variables described in this guide will have a substantial effect on system performance.

Remember that one system's `/etc/system` settings might not be applicable, either wholly or in part, to another system's environment. Carefully consider the values in the file with respect to the environment in which they will be applied. Make sure that you understand the behavior of a system before attempting to apply changes to the system variables that are described here.

We recommend that you start with an empty `/etc/system` file when moving to a new Solaris release. As a first step, add only those tunables that are required by in-house or third-party applications. After baseline testing has been established, evaluate system performance to determine if additional tunable settings are required.



**Caution** – The tunable parameters described in this book can and do change from Solaris release to Solaris release. Publication of these tunable parameters does not preclude changes to the tunable parameters and their descriptions without notice.

---

# Tuning Format of Tunable Parameters Descriptions

The format for the description of each tunable parameter is as follows:

- Parameter Name
- Description
- Data Type
- Default
- Range
- Units
- Dynamic?
- Validation
- Implicit
- When to Change
- Zone Configuration
- Commitment Level
- Change History

*Parameter Name* Is the exact name that is typed in the `/etc/system` file, or found in the `/etc/default/facility` file.

Most parameters names are of the form *parameter* where the parameter name does not contain a colon (:). These names refer to variables in the core portion of the kernel. If the name does contain a colon, the characters to the left of the colon reference the name of a loadable module. The name of the parameter within the module consists of the characters to the right of the colon. For example:

*module\_name:variable*

*Description* Briefly describes what the parameter does or controls.

*Data Type* Indicates the signed or unsigned short integer or long integer with the following distinctions:

- On a system that runs a 32-bit kernel, a long integer is the same size as an integer.
- On a system that runs a 64-bit kernel, a long integer is twice the width in bits as an integer. For example, an unsigned integer = 32 bits, an unsigned long integer = 64 bits.

*Units* (Optional) Describes the unit type.

*Default* What the system uses as the default value.

*Range* Specifies the possible range allowed by system validation or the bounds of the data type.

	<ul style="list-style-type: none"> <li>▪ MAXINT – A shorthand description for the maximum value of a signed integer (2,147,483,647)</li> <li>▪ MAXUINT – A shorthand description for the maximum value of an unsigned integer (4,294,967,295)</li> </ul>
Dynamic?	Yes, if the parameter can be changed on a running system with the mdb or kmdb debugger. No, if the parameter is a boot time initialization only.
Validation	Checks that the system applies to the value of the variable either as specified in the <code>/etc/system</code> file or the default value, as well as when the validation is applied.
Implicit	(Optional) Provides unstated constraints that might exist on the parameter, especially in relation to other parameters.
When to Change	Explains why someone might want to change this value. Includes error messages or return codes.
Zone Configuration	Identifies whether the parameter can be set in a exclusive-IP zone or must be set in the global zone. None of the parameters can be set in shared-IP zones.
Commitment Level	Identifies the stability of the interface. Many of the parameters in this manual are still evolving and are classified as unstable. For more information, see <a href="#">attributes(5)</a> .
Change History	(Optional) Contains a link to the Change History appendix, if applicable.

## Tuning the Solaris Kernel

The following table describes the different ways tunable parameters can be applied.

Apply Tunable Parameters in These Ways	For More Information
Modify the <code>/etc/system</code> file	<a href="#">“/etc/system File” on page 22</a>
Use the kernel debugger (kmdb)	<a href="#">“kmdb Command” on page 23</a>
Use the modular debugger (mdb)	<a href="#">“mdb Command” on page 23</a>
Use the <code>ipadm</code> command to set TCP/IP parameters	<a href="#">Chapter 4, “Internet Protocol Suite Tunable Parameters”</a>
Modify the <code>/etc/default</code> files	<a href="#">“Tuning NCA Parameters” on page 163</a>

## **/etc/system File**

The `/etc/system` file provides a static mechanism for adjusting the values of kernel parameters. Values specified in this file are read at boot time and are applied. Any changes that are made to the file are not applied to the operating system until the system is rebooted.

Prior to the Solaris 8 release, `/etc/system` entries that set the values of parameters were applied in two phases:

- The first phase obtains various bootstrap parameters (for example, `maxusers`) to initialize key system parameters.
- The second phase calculates the base configuration by using the bootstrap parameters, and all values specified in the `/etc/system` file are applied. In the case of the bootstrap parameters, reapplied values replace the values that are calculated or reset in the initialization phase.

The second phase sometimes caused confusion to users and administrators by setting parameters to values that seem to be impermissible or by assigning values to parameters (for example, `max_nprocs`) that have a value overridden during the initial configuration.

Starting in the Solaris 8 release, one pass is made to set all the values before the configuration parameters are calculated.

### **Example—Setting a Parameter in /etc/system**

The following `/etc/system` entry sets the ZFS ARC maximum (`zfs_arc_max`) to 30 GB.

```
set zfs:zfs_arc_max = 0x78000000
```

### **Recovering From an Incorrect Value**

Make a copy of the `/etc/system` file before modifying it so that you can easily recover from incorrect value. For example:

```
# cp /etc/system /etc/system.good
```

If a value specified in the `/etc/system` file causes the system to become unbootable, you can recover with the following command:

```
ok boot -a
```

This command causes the system to ask for the name of various files used in the boot process. Press the Return key to accept the default values until the name of the `/etc/system` file is requested. When the Name of system file `[/etc/system]:` prompt is displayed, type the name of the good `/etc/system` file or `/dev/null`:

```
Name of system file [/etc/system]: /etc/system.good
```

If `/dev/null` is specified, this path causes the system to attempt to read from `/dev/null` for its configuration information. Because this file is empty, the system uses the default values. After the system is booted, the `/etc/system` file can be corrected.

For more information on system recovery, see *System Administration Guide: Basic Administration*.

## kmdb Command

`kmdb` is a interactive kernel debugger with the same general syntax as `mdb`. An advantage of interactive kernel debugger is that you can set breakpoints. When a breakpoint is reached, you can examine data or step through the execution of kernel code.

`kmdb` can be loaded and unloaded on demand. You do not have to reboot the system to perform interactive kernel debugging, as was the case with `kadb`.

For more information, see `kmdb(1)`.

## mdb Command

Starting with the Solaris 8 release is the modular debugger, `mdb`, is unique among Solaris debuggers because it is easily extensible. A programming API is available that allows compilation of modules to perform desired tasks within the context of the debugger.

`mdb` also includes a number of desirable usability features, including command-line editing, command history, built-in output pager, syntax checking, and command pipelining. `mdb` is the recommended post-mortem debugger for the kernel.

For more information, see `mdb(1)`.

### Example—Using `mdb` to Change a Value

To change the value of the integer parameter `maxusers` from 495 to 512, do the following:

```
# mdb -kw
Loading modules: [ unix krtld genunix ip logindmux ptm nfs ipc lofs ]
> maxusers/D
maxusers:
maxusers:          495
> maxusers/W 200
maxusers:          0x1ef          =          0x200
> $q
```

Replace `maxusers` with the actual address of the item to be changed, as well as the value the parameter is to be set to.

For more information on using the modular debugger, see the *Solaris Modular Debugger Guide*.

When using either `kldb` or `mdb` debugger, the module name prefix is not required. After a module is loaded, its symbols form a common name space with the core kernel symbols and any other previously loaded module symbols.

For example, `ufs : ufs_WRITES` would be accessed as `ufs_WRITES` in each debugger (assuming the UFS module is loaded). The `ufs :` prefix is required when set in the `/etc/system` file.

## Special Solaris tune and var Structures

Solaris tunable parameters come in a variety of forms. The tune structure defined in the `/usr/include/sys/tuneable.h` file is the runtime representation of `tune_t_fsflushr`, `tune_t_minarmem`, and `tune_t_flkrec`. After the kernel is initialized, all references to these variables are found in the appropriate field of the tune structure.

Various documents (for example, previous versions of *Solaris System Administration Guide, Volume 2*) have stated that the proper way to set parameters in the tune structure is to use the syntax, `tune:field-name` where *field-name* is replaced by the actual parameter name listed above. This process silently fails. The proper way to set parameters for this structure at boot time is to initialize the special parameter that corresponds to the desired field name. The system initialization process then loads these values into the tune structure.

A second structure into which various tunable parameters are placed is the `var` structure named `v`. You can find the definition of a `var` structure in the `/usr/include/sys/var.h` file. The runtime representation of variables such as `autoup` and `bufhwm` is stored here.

Do not change either the `tune` or `v` structure on a running system. Changing any field in these structures on a running system might cause the system to panic.

## Viewing Solaris System Configuration Information

Several tools are available to examine system configuration information. Some tools require superuser privilege. Other tools can be run by a non-privileged user. Every structure and data item can be examined with the kernel debugger by using `mdb` on a running system or by booting under `kldb`.

For more information, see [mdb\(1\)](#) or [kadb\(1M\)](#).

### sysdef Command

The `sysdef` command provides the values of System V IPC settings, STREAMS tunables, process resource limits, and portions of the `tune` and `v` structures. For example, the `sysdef` “Tunable Parameters” section from on a 512-MB Sun Ultra 80 system is as follows:



---

334561280	maximum memory allowed in buffer cache (bufhwm)
30000	maximum number of processes (v.v_proc)
99	maximum global priority in sys class (MAXCLSPRI)
29995	maximum processes per user id (v.v_maxup)
30	auto update time limit in seconds (NAUTOUP)
25	page stealing low water mark (GPGSLO)
1	fsflush run rate (FSFLUSHR)
25	minimum resident memory for avoiding deadlock (MINARMEM)
25	minimum swapable memory for avoiding deadlock (MINASMEM)

For more information, see [sysdef\(1M\)](#).

## kstat Utility

kstats are data structures maintained by various kernel subsystems and drivers. They provide a mechanism for exporting data from the kernel to user programs without requiring that the program read kernel memory or have superuser privilege. For more information, see [kstat\(1M\)](#) or [kstat\(3KSTAT\)](#).

---

**Note** – kstat data structures with `system_pages` name in the `unix` module do not report statistics for `cachefree`. `cachefree` is not supported, starting in the Solaris 9 release.

---



# Oracle Solaris Kernel Tunable Parameters

---

This chapter describes most of the Oracle Solaris kernel tunable parameters.

- “General Kernel and Memory Parameters” on page 28
- “fsflush and Related Parameters” on page 34
- “Process-Sizing Parameters” on page 38
- “Paging-Related Parameters” on page 42
- “Swapping-Related Parameters” on page 54
- “Kernel Memory Allocator” on page 55
- “General Driver Parameters” on page 58
- “Network Driver Parameters” on page 60
- “General I/O Parameters” on page 65
- “General File System Parameters” on page 67
- “UFS Parameters” on page 71
- “TMPFS Parameters” on page 78
- “Pseudo Terminals” on page 79
- “STREAMS Parameters” on page 82
- “System V Message Queues” on page 83
- “System V Semaphores” on page 83
- “System V Shared Memory” on page 84
- “Scheduling” on page 85
- “Timers” on page 85
- “sun4u or sun4v Specific Parameters” on page 86
- “Locality Group Parameters” on page 90

## Where to Find Tunable Parameter Information

Tunable Parameter	For Information
NFS tunable parameters	<a href="#">Chapter 3, “NFS Tunable Parameters”</a>
Internet Protocol Suite tunable parameters	<a href="#">Chapter 4, “Internet Protocol Suite Tunable Parameters”</a>
Network Cache and Accelerator (NCA) tunable parameters	<a href="#">Chapter 5, “Network Cache and Accelerator Tunable Parameters”</a>

## General Kernel and Memory Parameters

This section describes general kernel parameters that are related to physical memory and stack configuration.

### physmem

Description	Modifies the system's configuration of the number of physical pages of memory after the Solaris OS and firmware are accounted for.
Data Type	Unsigned long
Default	Number of usable pages of physical memory available on the system, not counting the memory where the core kernel and data are stored
Range	1 to amount of physical memory on system
Units	Pages
Dynamic?	No
Validation	None
When to Change	Whenever you want to test the effect of running the system with less physical memory. Because this parameter does <i>not</i> take into account the memory used by the core kernel and data, as well as various other data structures allocated early in the startup process, the value of <code>physmem</code> should be less than the actual number of pages that represent the smaller amount of memory.
Commitment Level	Unstable

## zfs\_arc\_min

Description	Determines the minimum size of the ZFS Adjustable Replacement Cache (ARC). See also <a href="#">“zfs_arc_max” on page 29</a> .
Data Type	Unsigned Integer (64-bit)
Default	1/32nd of physical memory or 64 MB, whichever value is larger.
Range	64 MB to <code>zfs_arc_max</code>
Units	Bytes
Dynamic?	No
Validation	Yes, the range is validated.
When to Change	When a system's workload demand for memory fluctuates, the ZFS ARC caches data at a period of weak demand and then shrinks at a period of strong demand. However, ZFS does not shrink below the value of <code>zfs_arc_min</code> . The default value of <code>zfs_arc_min</code> is 12% of memory on large memory systems and so, can be a significant amount of memory. If a workload's highest memory usage requires more than 88% of system memory, consider tuning this parameter.
Commitment Level	Unstable
Change History	For information, see <a href="#">“zfs_arc_min (Oracle Solaris 11 Express)” on page 176</a> .

## zfs\_arc\_max

Description	Determines the maximum size of the ZFS Adjustable Replacement Cache (ARC). See also <a href="#">“zfs_arc_min” on page 29</a> .
Data Type	Unsigned Integer (64-bit)
Default	Three-fourths of memory on systems with less than 4 GB of memory physmem minus 1 GB on systems with greater than 4 GB of memory
Range	64 MB to <code>physmem</code>
Units	Bytes
Dynamic?	No
Validation	Yes, the range is validated.
When to Change	If a future memory requirement is significantly large and well defined, you might consider reducing the value of this parameter to cap the

ARC so that it does not compete with the memory requirement. For example, if you know that a future workload requires 20% of memory, it makes sense to cap the ARC such that it does not consume more than the remaining 80% of memory.

Commitment Level	Unstable
Change History	For information, see “ <a href="#">zfs_arc_max (Oracle Solaris 11 Express)</a> ” on page 176.

## default\_stksize

Description	Specifies the default stack size of all threads. No thread can be created with a stack size smaller than <code>default_stksize</code> . If <code>default_stksize</code> is set, it overrides <code>lwp_default_stksize</code> . See also “ <a href="#">lwp_default_stksize</a> ” on page 31.
Data Type	Integer
Default	<ul style="list-style-type: none"><li>▪ 3 x PAGESIZE on SPARC systems</li><li>▪ 2 x PAGESIZE on x86 systems</li><li>▪ 5 x PAGESIZE on AMD64 systems</li></ul>
Range	Minimum is the default values: <ul style="list-style-type: none"><li>▪ 3 x PAGESIZE on SPARC systems</li><li>▪ 2 x PAGESIZE on x86 systems</li><li>▪ 5 x PAGESIZE on AMD64 systems</li></ul> Maximum is 32 times the default value.
Units	Bytes in multiples of the value returned by the <code>getpagesize</code> parameter. For more information, see <a href="#">getpagesize(3C)</a> .
Dynamic?	Yes. Affects threads created after the variable is changed.
Validation	Must be greater than or equal to 8192 and less than or equal to 262,144 (256 x 1024). Also must be a multiple of the system page size. If these conditions are not met, the following message is displayed: <pre>Illegal stack size, Using N</pre> The value of <i>N</i> is the default value of <code>default_stksize</code> .
When to Change	When the system panics because it has run out of stack space. The best solution for this problem is to determine why the system is running out of space and then make a correction.

Increasing the default stack size means that almost every kernel thread will have a larger stack, resulting in increased kernel memory consumption for no good reason. Generally, that space will be unused. The increased consumption means other resources that are competing for the same pool of memory will have the amount of space available to them reduced, possibly decreasing the system's ability to perform work. Among the side effects is a reduction in the number of threads that the kernel can create. This solution should be treated as no more than an interim workaround until the root cause is remedied.

Commitment Level    Unstable

## **lwp\_default\_stksize**

Description	Specifies the default value of the stack size to be used when a kernel thread is created, and when the calling routine does not provide an explicit size to be used.
Data Type	Integer
Default	<ul style="list-style-type: none"> <li>▪ 8192 for x86 platforms</li> <li>▪ 24,576 for SPARC platforms</li> <li>▪ 20,480 for AMD64 platforms</li> </ul>
Range	<p>Minimum is the default values:</p> <ul style="list-style-type: none"> <li>▪ 3 x PAGESIZE on SPARC systems</li> <li>▪ 2 x PAGESIZE on x86 systems</li> <li>▪ 5 x PAGESIZE on AMD64 systems</li> </ul> <p>Maximum is 32 times the default value.</p>
Units	Bytes in multiples of the value returned by the <code>getpagesize</code> parameter. For more information, see <a href="#">getpagesize(3C)</a> .
Dynamic?	Yes. Affects threads created after the variable is changed.
Validation	<p>Must be greater than or equal to 8192 and less than or equal to 262,144 (256 x 1024). Also must be a multiple of the system page size. If these conditions are not met, the following message is displayed:</p> <pre>Illegal stack size, Using N</pre> <p>The value of <i>N</i> is the default value of <code>lwp_default_stksize</code>.</p>

When to Change	<p>When the system panics because it has run out of stack space. The best solution for this problem is to determine why the system is running out of space and then make a correction.</p> <p>Increasing the default stack size means that almost every kernel thread will have a larger stack, resulting in increased kernel memory consumption for no good reason. Generally, that space will be unused. The increased consumption means other resources that are competing for the same pool of memory will have the amount of space available to them reduced, possibly decreasing the system's ability to perform work. Among the side effects is a reduction in the number of threads that the kernel can create. This solution should be treated as no more than an interim workaround until the root cause is remedied.</p>
Commitment Level	Unstable

## logevent\_max\_q\_sz

Description	Maximum number of system events allowed to be queued and waiting for delivery to the syseventd daemon. Once the size of the system event queue reaches this limit, no other system events are allowed on the queue.
Data Type	Integer
Default	5000
Range	0 to MAXINT
Units	System events
Dynamic?	Yes
Validation	<p>The system event framework checks this value every time a system event is generated by <code>ddi_log_sysevent</code> and <code>sysevent_post_event</code>.</p> <p>For more information, see <a href="#">ddi_log_sysevent(9F)</a> and <a href="#">sysevent_post_event(3SYSEVENT)</a>.</p>
When to Change	When error log messages indicate that a system event failed to be logged, generated, or posted.
Commitment Level	Unstable



## segkpsize

Description	Specifies the amount of kernel pageable memory available. This memory is used primarily for kernel thread stacks. Increasing this number allows either larger stacks for the same number of threads or more threads. This parameter can only be set on a system running a 64-bit kernel. A system running a 64-bit kernel uses a default stack size of 24 KB.
Data Type	Unsigned long
Default	64-bit kernels, 2 GB 32-bit kernels, 512 MB
Range	64-bit kernels, 512 MB to 24 GB
Units	8-KB pages
Dynamic?	No
Validation	Value is compared to minimum and maximum sizes (512 MB and 24 GB for 64-bit systems). If smaller than the minimum or larger than the maximum, it is reset to 2 GB. A message to that effect is displayed.  The actual size used in creation of the cache is the lesser of the value specified in <code>segkpsize</code> after the validation checking or 50 percent of physical memory.
When to Change	Required to support large numbers of processes on a system. The default size of 2 GB, assuming at least 1 GB of physical memory is present. This default size allows creation of 24-KB stacks for more than 87,000 kernel threads. The size of a stack in a 64-bit kernel is the same, whether the process is a 32-bit process or a 64-bit process. If more than this number is needed, <code>segkpsize</code> can be increased, assuming sufficient physical memory exists.
Commitment Level	Unstable

## noexec\_user\_stack

Description	Enables the stack to be marked as nonexecutable, which helps make buffer-overflow attacks more difficult.
-------------	---

A Solaris system running a 64-bit kernel makes the stacks of all 64-bit applications nonexecutable by default. Setting this parameter is necessary to make 32-bit applications nonexecutable on systems running 64-bit or 32-bit kernels.

---

**Note** – This parameter exists on all systems running the Solaris 2.6, 7, 8, 9, or 10 releases, but it is only effective on 64-bit SPARC and AMD64 architectures.

---

Data Type	Signed integer
Default	0 (disabled)
Range	0 (disabled) or 1 (enabled)
Units	Toggle (on/off)
Dynamic?	Yes. Does not affect currently running processes, only processes created after the value is set.
Validation	None
When to Change	Should be enabled at all times unless applications are deliberately placing executable code on the stack without using <code>mprotect</code> to make the stack executable. For more information, see <a href="#">mprotect(2)</a> .
Commitment Level	Unstable

## fsflush and Related Parameters

This section describes `fsflush` and related tunables.

### fsflush

The system daemon, `fsflush`, runs periodically to do three main tasks:

1. On every invocation, `fsflush` flushes dirty file system pages over a certain age to disk.
2. On every invocation, `fsflush` examines a portion of memory and causes modified pages to be written to their backing store. Pages are written if they are modified and if they do not meet one of the following conditions:
  - Pages are kernel page
  - Pages are free
  - Pages are locked

- Pages are associated with a swap device
- Pages are currently involved in an I/O operation

The net effect is to flush pages from files that are mapped with `mmap` with write permission and that have actually been changed.

Pages are flushed to backing store but left attached to the process using them. This will simplify page reclamation when the system runs low on memory by avoiding delay for writing the page to backing store before claiming it, if the page has not been modified since the flush.

3. `fsflush` writes file system metadata to disk. This write is done every  $n$ th invocation, where  $n$  is computed from various configuration variables. See “`tune_t_fsflushr`” on page 35 and “`autoup`” on page 36 for details.

The following features are configurable:

- Frequency of invocation (`tune_t_fsflushr`)
- Whether memory scanning is executed (`dopageflush`)
- Whether file system data flushing occurs (`doiflush`)
- The frequency with which file system data flushing occurs (`autoup`)

For most systems, memory scanning and file system metadata synchronizing are the dominant activities for `fsflush`. Depending on system usage, memory scanning can be of little use or consume too much CPU time.

## tune\_t\_fsflushr

Description	Specifies the number of seconds between <code>fsflush</code> invocations
Data Type	Signed integer
Default	1
Range	1 to MAXINT
Units	Seconds
Dynamic?	No
Validation	If the value is less than or equal to zero, the value is reset to 1 and a warning message is displayed. This check is done only at boot time.
When to Change	See the <code>autoup</code> parameter.
Commitment Level	Unstable

## autoup

Description	<p>Along with <code>tune_t_flushr</code>, <code>autoup</code> controls the amount of memory examined for dirty pages in each invocation and frequency of file system synchronizing operations.</p> <p>The value of <code>autoup</code> is also used to control whether a buffer is written out from the free list. Buffers marked with the <code>B_DELWRI</code> flag (which identifies file content pages that have changed) are written out whenever the buffer has been on the list for longer than <code>autoup</code> seconds. Increasing the value of <code>autoup</code> keeps the buffers in memory for a longer time.</p>
Data Type	Signed integer
Default	30
Range	1 to MAXINT
Units	Seconds
Dynamic?	No
Validation	If <code>autoup</code> is less than or equal to zero, it is reset to 30 and a warning message is displayed. This check is done only at boot time.
Implicit	<p><code>autoup</code> should be an integer multiple of <code>tune_t_fsflushr</code>. At a minimum, <code>autoup</code> should be at least 6 times the value of <code>tune_t_fsflushr</code>. If not, excessive amounts of memory are scanned each time <code>fsflush</code> is invoked.</p> <p>The total system pages multiplied by <code>tune_t_fsflushr</code> should be greater than or equal to <code>autoup</code> to cause memory to be checked if <code>dopageflush</code> is non-zero.</p>
When to Change	<p>Here are several potential situations for changing <code>autoup</code>, <code>tune_t_fsflushr</code>, or both:</p> <ul style="list-style-type: none"><li>▪ Systems with large amounts of memory – In this case, increasing <code>autoup</code> reduces the amount of memory scanned in each invocation of <code>fsflush</code>.</li><li>▪ Systems with minimal memory demand – Increasing both <code>autoup</code> and <code>tune_t_fsflushr</code> reduces the number of scans made. <code>autoup</code> should be increased also to maintain the current ratio of <code>autoup</code> / <code>tune_t_fsflushr</code>.</li></ul>

- Systems with large numbers of transient files (for example, mail servers or software build machines) – If large numbers of files are created and then deleted, fsflush might unnecessarily write data pages for those files to disk.

Commitment Level Unstable

## dopageflush

Description	Controls whether memory is examined for modified pages during fsflush invocations. In each invocation of fsflush, the number of physical memory pages in the system is determined. This number might have changed because of a dynamic reconfiguration operation. Each invocation scans by using this algorithm: total number of pages $\times$ $\text{tune\_t\_fsflushr} / \text{autoup}$ pages
Data Type	Signed integer
Default	1 (enabled)
Range	0 (disabled) or 1 (enabled)
Units	Toggle (on/off)
Dynamic?	Yes
Validation	None
When to Change	If the system page scanner rarely runs, which is indicated by a value of 0 in the sr column of vmsstat output.
Commitment Level	Unstable
Change History	For information, see <a href="#">“dopageflush (All Solaris Releases)”</a> on page 176.

## doiflush

Description	Controls whether file system metadata syncs will be executed during fsflush invocations. This synchronization is done every $N$ th invocation of fsflush where $N = (\text{autoup} / \text{tune\_t\_fsflushr})$ . Because this algorithm is integer division, if $\text{tune\_t\_fsflushr}$ is greater than $\text{autoup}$ , a synchronization is done on every invocation of fsflush because the code checks to see if its iteration counter is greater than or equal to $N$ . Note that $N$ is computed once on invocation of fsflush. Later changes to $\text{tune\_t\_fsflushr}$ or $\text{autoup}$ have no effect on the frequency of synchronization operations.
-------------	---

Data Type	Signed integer
Default	1 (enabled)
Range	0 (disabled) or 1 (enabled)
Units	Toggle (on/off)
Dynamic?	Yes
Validation	None
When to Change	<p>When files are frequently modified over a period of time and the load caused by the flushing perturbs system behavior.</p> <p>Files whose existence, and therefore consistency of state, does not matter if the system reboots are better kept in a TMPFS file system (for example, /tmp). Inode traffic can be reduced on systems, starting in the Solaris 7 release, by using the <code>mount -noatime</code> option. This option eliminates inode updates when the file is accessed.</p> <p>For a system engaged in realtime processing, you might want to disable this option and use explicit application file synchronizing to achieve consistency.</p>
Commitment Level	Unstable

## Process-Sizing Parameters

Several parameters (or variables) are used to control the number of processes that are available on the system and the number of processes that an individual user can create. The foundation parameter is `maxusers`. This parameter drives the values assigned to `max_nprocs` and `maxuprc`.

### **maxusers**

Description	<p>Originally, <code>maxusers</code> defined the number of logged in users the system could support. When a kernel was generated, various tables were sized based on this setting. Current Solaris releases do much of its sizing based on the amount of memory on the system. Thus, much of the past use of <code>maxusers</code> has changed. A number of subsystems that are still derived from <code>maxusers</code>:</p> <ul style="list-style-type: none"><li>▪ The maximum number of processes on the system</li><li>▪ The number of quota structures held in the system</li><li>▪ The size of the directory name look-up cache (DNLC)</li></ul>
-------------	---

Data Type	Signed integer
Default	Lesser of the amount of memory in MB or 2048
Range	1 to 2048, based on physical memory if not set in the <code>/etc/system</code> file 1 to 4096, if set in the <code>/etc/system</code> file
Units	Users
Dynamic?	No. After computation of dependent parameters is done, <code>maxusers</code> is never referenced again.
Validation	None
When to Change	When the default number of user processes derived by the system is too low. This situation is evident when the following message displays on the system console:  out of processes  You might also change this parameter when the default number of processes is too high, as in these situations: <ul style="list-style-type: none"> <li>▪ Database servers that have a lot of memory and relatively few running processes can save system memory when the default value of <code>maxusers</code> is reduced.</li> <li>▪ If file servers have a lot of memory and few running processes, you might reduce this value. However, you should explicitly set the size of the DNLC. See “<a href="#">ncsize</a>” on page 67.</li> <li>▪ If compute servers have a lot of memory and few running processes, you might reduce this value.</li> </ul>
Commitment Level	Unstable

## reserved\_procs

Description	Specifies the number of system process slots to be reserved in the process table for processes with a UID of root (0). For example, <code>fsflush</code> has a UID of root (0).
Data Type	Signed integer
Default	5
Range	5 to MAXINT
Units	Processes
Dynamic?	No. Not used after the initial parameter computation.

Validation	Starting in the Solaris 8 release, any <code>/etc/system</code> setting is honored.
Commitment Level	Unstable
When to Change	Consider increasing to 10 + the normal number of UID 0 (root) processes on system. This setting provides some cushion should it be necessary to obtain a root shell when the system is otherwise unable to create user-level processes.

## **pidmax**

Description	<p>Specifies the value of the largest possible process ID. Valid for Solaris 8 and later releases.</p> <p><code>pidmax</code> sets the value for the <code>maxpid</code> variable. Once <code>maxpid</code> is set, <code>pidmax</code> is ignored. <code>maxpid</code> is used elsewhere in the kernel to determine the maximum process ID and for validation checking.</p> <p>Any attempts to set <code>maxpid</code> by adding an entry to the <code>/etc/system</code> file have no effect.</p>
Data Type	Signed integer
Default	30,000
Range	266 to 999,999
Units	Processes
Dynamic?	No. Used only at boot time to set the value of <code>pidmax</code> .
Validation	Yes. Value is compared to the value of <code>reserved_procs</code> and 999,999. If less than <code>reserved_procs</code> or greater than 999,999, the value is set to 999,999.
Implicit	<code>max_nprocs</code> range checking ensures that <code>max_nprocs</code> is always less than or equal to this value.
When to Change	Required to enable support for more than 30,000 processes on a system.
Commitment Level	Unstable



## max\_nprocs

Description	<p>Specifies the maximum number of processes that can be created on a system. Includes system processes and user processes. Any value specified in <code>/etc/system</code> is used in the computation of <code>maxuprc</code>.</p> <p>This value is also used in determining the size of several other system data structures. Other data structures where this parameter plays a role are as follows:</p> <ul style="list-style-type: none"> <li>▪ Determining the size of the directory name lookup cache (if <code>ncsize</code> is not specified)</li> <li>▪ Allocating disk quota structures for UFS (if <code>ndquot</code> is not specified)</li> <li>▪ Verifying that the amount of memory used by configured system V semaphores does not exceed system limits</li> <li>▪ Configuring Hardware Address Translation resources for x86 platforms.</li> </ul>
Data Type	Signed integer
Default	$10 + (16 \times \text{maxusers})$
Range	266 to value of <code>maxpid</code>
Dynamic?	No
Validation	Yes. The value is compared to <code>maxpid</code> and set to <code>maxpid</code> if it is larger. On x86 platforms, an additional check is made against a platform-specific value. <code>max_nprocs</code> is set to the smallest value in the triplet ( <code>max_nprocs</code> , <code>maxpid</code> , platform value). Both SPARC and x86 platforms use 65,534 as the platform value.
When to Change	Changing this parameter is one of the steps necessary to enable support for more than 30,000 processes on a system.
Commitment Level	Unstable

## maxuprc

Description	Specifies the maximum number of processes that can be created on a system by any one user.
Data Type	Signed integer
Default	<code>max_nprocs - reserved_procs</code>
Range	1 to <code>max_nprocs - reserved_procs</code>

Units	Processes
Dynamic?	No
Validation	Yes. This value is compared to <code>max_nprocs - reserved_procs</code> and set to the smaller of the two values.
When to Change	When you want to specify a hard limit for the number of processes a user can create that is less than the default value of however many processes the system can create. Attempting to exceed this limit generates the following warning messages on the console or in the messages file:  <code>out of per-user processes for uid N</code>
Commitment Level	Unstable

## **ngroups\_max**

Description	Specifies the maximum number of supplemental groups per process.
Data Type	Signed integer
Default	16
Range	0 to 1024
Units	Groups
Dynamic?	No
Validation	No
When to Change	When you want to increase the maximum number of groups.  Keep in mind that if a particular user is assigned to more than 16 groups, the user might experience problems with <code>AUTH_SYS</code> credentials in an NFS environment.
Commitment Level	Unstable

## **Paging-Related Parameters**

The Solaris OS uses a demand paged virtual memory system. As the system runs, pages are brought into memory as needed. When memory becomes occupied above a certain threshold and demand for memory continues, paging begins. Paging goes through several levels that are controlled by certain parameters.

The general paging algorithm is as follows:

- A memory deficit is noticed. The page scanner thread runs and begins to walk through memory. A two-step algorithm is employed:
  1. A page is marked as unused.
  2. If still unused after a time interval, the page is viewed as a subject for reclaim.

If the page has been modified, a request is made to the pageout thread to schedule the page for I/O. Also, the page scanner continues looking at memory. Pageout causes the page to be written to the page's backing store and placed on the free list. When the page scanner scans memory, no distinction is made as to the origin of the page. The page might have come from a data file, or it might represent a page from an executable's text, data, or stack.

- As memory pressure on the system increases, the algorithm becomes more aggressive in the pages it will consider as candidates for reclamation and in how frequently the paging algorithm runs. (For more information, see “[fastscan](#)” on page 50 and “[slowscan](#)” on page 50.) As available memory falls between the range `lotsfree` and `minfree`, the system linearly increases the amount of memory scanned in each invocation of the pageout thread from the value specified by `slowscan` to the value specified by `fastscan`. The system uses the `desfree` parameter to control a number of decisions about resource usage and behavior.

The system initially constrains itself to use no more than 4 percent of one CPU for pageout operations. As memory pressure increases, the amount of CPU time consumed in support of pageout operations linearly increases until a maximum of 80 percent of one CPU is consumed. The algorithm looks through some amount of memory between `slowscan` and `fastscan`, then stops when one of the following occurs:

- Enough pages have been found to satisfy the memory shortfall.
- The planned number of pages have been looked at.
- Too much time has elapsed.

If a memory shortfall is still present when pageout finishes its scan, another scan is scheduled for 1/4 second in the future.

The configuration mechanism of the paging subsystem was changed, starting in the Solaris 9 release. Instead of depending on a set of predefined values for `fastscan`, `slowscan`, and `handspreadpages`, the system determines the appropriate settings for these parameters at boot time. Setting any of these parameters in the `/etc/system` file can cause the system to use less than optimal values.




---

**Caution** – Remove all tuning of the VM system from the `/etc/system` file. Run with the default settings and determine if it is necessary to adjust any of these parameters. Do not set either `cachefree` or `priority_paging`. They have been removed, starting in the Solaris 9 release.

---

Beginning in the Solaris 7 5/99 release, dynamic reconfiguration (DR) for CPU and memory is supported. A system in a DR operation that involves the addition or deletion of memory recalculates values for the relevant parameters, unless the parameter has been explicitly set in `/etc/system`. In that case, the value specified in `/etc/system` is used, unless a constraint on the value of the variable has been violated. In this case, the value is reset.

## lotsfree

Description	Serves as the initial trigger for system paging to begin. When this threshold is crossed, the page scanner wakes up to begin looking for memory pages to reclaim.
Data Type	Unsigned long
Default	The greater of 1/64th of physical memory or 512 KB
Range	<p>The minimum value is 512 KB or 1/64th of physical memory, whichever is greater, expressed as pages using the page size returned by <code>getpagesize(3C)</code>.</p> <p>The maximum value is the number of physical memory pages. The maximum value should be no more than 30 percent of physical memory. The system does not enforce this range, other than that described in the Validation section.</p>
Units	Pages
Dynamic?	Yes, but dynamic changes are lost if a memory-based DR operation occurs.
Validation	If <code>lotsfree</code> is greater than the amount of physical memory, the value is reset to the default.
Implicit	The relationship of <code>lotsfree</code> being greater than <code>desfree</code> , which is greater than <code>minfree</code> , should be maintained at all times.
When to Change	<p>When demand for pages is subject to sudden sharp spikes, the memory algorithm might be unable to keep up with demand. One workaround is to start reclaiming memory at an earlier time. This solution gives the paging system some additional margin.</p> <p>A rule of thumb is to set this parameter to 2 times what the system needs to allocate in a few seconds. This parameter is workload dependent. A DBMS server can probably work fine with the default settings. However, you might need to adjust this parameter for a system doing heavy file system I/O.</p>

	For systems with relatively static workloads and large amounts of memory, lower this value. The minimum acceptable value is 512 KB, expressed as pages using the page size returned by <code>getpagesize</code> .
Commitment Level	Unstable

## desfree

Description	Specifies the preferred amount of memory to be free at all times on the system.
Data Type	Unsigned integer
Default	<code>lotsfree / 2</code>
Range	The minimum value is 256 KB or 1/128th of physical memory, whichever is greater, expressed as pages using the page size returned by <code>getpagesize</code> .  The maximum value is the number of physical memory pages. The maximum value should be no more than 15 percent of physical memory. The system does not enforce this range other than that described in the Validation section.
Units	Pages
Dynamic?	Yes, unless dynamic reconfiguration operations that add or delete memory occur. At that point, the value is reset to the value provided in the <code>/etc/system</code> file or calculated from the new physical memory value.
Validation	If <code>desfree</code> is greater than <code>lotsfree</code> , <code>desfree</code> is set to <code>lotsfree / 2</code> . No message is displayed.
Implicit	The relationship of <code>lotsfree</code> being greater than <code>desfree</code> , which is greater than <code>minfree</code> , should be maintained at all times.
Side Effects	Several side effects can arise from increasing the value of this parameter. When the new value nears or exceeds the amount of available memory on the system, the following can occur: <ul style="list-style-type: none"> <li>▪ Asynchronous I/O requests are not processed, unless available memory exceeds <code>desfree</code>. Increasing the value of <code>desfree</code> can result in rejection of requests that otherwise would succeed.</li> <li>▪ NFS asynchronous writes are executed as synchronous writes.</li> <li>▪ The swapper is awakened earlier, and the behavior of the swapper is biased towards more aggressive actions.</li> </ul>

- The system might not prefault as many executable pages into the system. This side effect results in applications potentially running slower than they otherwise would.

When to Change	For systems with relatively static workloads and large amounts of memory, lower this value. The minimum acceptable value is 256 KB, expressed as pages using the page size returned by <code>getpagesize</code> .
Commitment Level	Unstable

## **minfree**

Description	Specifies the minimum acceptable memory level. When memory drops below this number, the system biases allocations toward allocations necessary to successfully complete pageout operations or to swap processes completely out of memory. Either allocation denies or blocks other allocation requests.
Data Type	Unsigned integer
Default	<code>desfree / 2</code>
Range	<p>The minimum value is 128 KB or 1/256th of physical memory, whichever is greater, expressed as pages using the page size returned by <code>getpagesize</code>.</p> <p>The maximum value is the number of physical memory pages. The maximum value should be no more than 7.5 percent of physical memory. The system does not enforce this range other than that described in the Validation section.</p>
Units	Pages
Dynamic?	Yes, unless dynamic reconfiguration operations that add or delete memory occur. At that point, the value is reset to the value provided in the <code>/etc/system</code> file or calculated from the new physical memory value.
Validation	If <code>minfree</code> is greater than <code>desfree</code> , <code>minfree</code> is set to <code>desfree / 2</code> . No message is displayed.
Implicit	The relationship of <code>lotsfree</code> being greater than <code>desfree</code> , which is greater than <code>minfree</code> , should be maintained at all times.
When to Change	The default value is generally adequate. For systems with relatively static workloads and large amounts of memory, lower this value. The minimum acceptable value is 128 KB, expressed as pages using the page size returned by <code>getpagesize</code> .

Commitment Level Unstable

## throttlefree

Description	Specifies the memory level at which blocking memory allocation requests are put to sleep, even if the memory is sufficient to satisfy the request.
Data Type	Unsigned integer
Default	<code>minfree</code>
Range	The minimum value is 128 KB or 1/256th of physical memory, whichever is greater, expressed as pages using the page size returned by <code>getpagesize</code> .  The maximum value is the number of physical memory pages. The maximum value should be no more than 4 percent of physical memory. The system does not enforce this range other than that described in the Validation section.
Units	Pages
Dynamic?	Yes, unless dynamic reconfiguration operations that add or delete memory occur. At that point, the value is reset to the value provided in the <code>/etc/system</code> file or calculated from the new physical memory value.
Validation	If <code>throttlefree</code> is greater than <code>desfree</code> , <code>throttlefree</code> is set to <code>minfree</code> . No message is displayed.
Implicit	The relationship of <code>lotsfree</code> is greater than <code>desfree</code> , which is greater than <code>minfree</code> , should be maintained at all times.
When to Change	The default value is generally adequate. For systems with relatively static workloads and large amounts of memory, lower this value. The minimum acceptable value is 128 KB, expressed as pages using the page size returned by <code>getpagesize</code> . For more information, see <a href="#">getpagesize(3C)</a> .
Commitment Level	Unstable

## pageout\_reserve

Description	Specifies the number of pages reserved for the exclusive use of the pageout or scheduler threads. When available memory is less than this
-------------	---

	<p>value, nonblocking allocations are denied for any processes other than pageout or the scheduler. Pageout needs to have a small pool of memory for its use so it can allocate the data structures necessary to do the I/O for writing a page to its backing store. This variable was introduced in the Solaris 2.6 release to ensure that the system would be able to perform a pageout operation in the face of the most severe memory shortage.</p>
Data Type	Unsigned integer
Default	<code>throttlefree / 2</code>
Range	<p>The minimum value is 64 KB or 1/512th of physical memory, whichever is greater, expressed as pages using the page size returned by <code>getpagesize(3C)</code>.</p> <p>The maximum is the number of physical memory pages. The maximum value should be no more than 2 percent of physical memory. The system does not enforce this range, other than that described in the Validation section.</p>
Units	Pages
Dynamic?	Yes, unless dynamic reconfiguration operations that add or delete memory occur. At that point, the value is reset to the value provided in the <code>/etc/system</code> file or calculated from the new physical memory value.
Validation	If <code>pageout_reserve</code> is greater than <code>throttlefree / 2</code> , <code>pageout_reserve</code> is set to <code>throttlefree / 2</code> . No message is displayed.
Implicit	The relationship of <code>lotsfree</code> being greater than <code>desfree</code> , which is greater than <code>minfree</code> , should be maintained at all times.
When to Change	The default value is generally adequate. For systems with relatively static workloads and large amounts of memory, lower this value. The minimum acceptable value is 64 KB, expressed as pages using the page size returned by <code>getpagesize</code> .
Commitment Level	Unstable

## **pages\_pp\_maximum**

Description	Defines the number of pages that must be unlocked. If a request to lock pages would force available memory below this value, that request is refused.
Data Type	Unsigned long



Default	The greater of ( <code>tune_t_minarmem + 100</code> and [4% of memory available at boot time + 4 MB])
Range	Minimum value enforced by the system is <code>tune_t_minarmem + 100</code> . The system does not enforce a maximum value.
Units	Pages
Dynamic?	Yes, unless dynamic reconfiguration operations that add or delete memory occur. At that point, the value is reset to the value provided in the <code>/etc/system</code> file or was calculated from the new physical memory value.
Validation	If the value specified in the <code>/etc/system</code> file or the calculated default is less than <code>tune_t_minarmem + 100</code> , the value is reset to <code>tune_t_minarmem + 100</code> .  No message is displayed if the value from the <code>/etc/system</code> file is increased. Validation is done only at boot time and during dynamic reconfiguration operations that involve adding or deleting memory.
When to Change	When memory-locking requests fail or when attaching to a shared memory segment with the <code>SHARE_MMU</code> flag fails, yet the amount of memory available seems to be sufficient.  Excessively large values can cause memory locking requests ( <code>mlock</code> , <code>mlockall</code> , and <code>mlockntl</code> ) to fail unnecessarily. For more information, see <a href="#">mlock(3C)</a> , <a href="#">mlockall(3C)</a> , and <a href="#">mlockntl(2)</a> .
Commitment Level	Unstable

## **tune\_t\_minarmem**

Description	Defines the minimum available resident (not swappable) memory to maintain necessary to avoid deadlock. Used to reserve a portion of memory for use by the core of the OS. Pages restricted in this way are not seen when the OS determines the maximum amount of memory available.
Data Type	Signed integer
Default	25
Range	1 to physical memory
Units	Pages
Dynamic?	No

Validation	None. Large values result in wasted physical memory.
When to Change	The default value is generally adequate. Consider increasing the default value if the system locks up and debugging information indicates that no memory was available.
Commitment Level	Unstable

## **fastscan**

Description	Defines the maximum number of pages per second that the system looks at when memory pressure is highest.
Data Type	Signed integer
Default	After the system is booted, fastscan is set to 64 MB. Then this value is automatically reset to the number of pages that the scanner can scan in one second by using 10% of a CPU. If this derived value is more than half the system's physical memory, the default value is limited to half the system's physical memory.
Range	64 MB to half the system's physical memory
Units	Pages
Dynamic?	Yes, unless dynamic reconfiguration operations that add or delete memory occur. At that point, the value is reset to the value provided by <code>/etc/system</code> or calculated from the new physical memory value.
Validation	The maximum value is the lesser of 64 MB and 1/2 of physical memory.
When to Change	When more aggressive scanning of memory is preferred during periods of memory shortfall, especially when the system is subject to periods of intense memory demand or when performing heavy file I/O.
Commitment Level	Unstable

## **slowscan**

Description	Defines the minimum number of pages per second that the system looks at when attempting to reclaim memory.
Data Type	Signed integer
Default	The smaller of 1/20th of physical memory in pages and 100.
Range	1 to <code>fastscan / 2</code>

Units	Pages
Dynamic?	Yes, unless dynamic reconfiguration operations that add or delete memory occur. At that point, the value is reset to the value provided in the <code>/etc/system</code> file or calculated from the new physical memory value.
Validation	If <code>slowscan</code> is larger than <code>fastscan / 2</code> , <code>slowscan</code> is reset to <code>fastscan / 2</code> . No message is displayed.
When to Change	When more aggressive scanning of memory is preferred during periods of memory shortfall, especially when the system is subject to periods of intense memory demand.
Commitment Level	Unstable

## **min\_percent\_cpu**

Description	Defines the minimum percentage of CPU that pageout can consume. This parameter is used as the starting point for determining the maximum amount of time that can be consumed by the page scanner.
Data Type	Signed integer
Default	4
Range	1 to 80
Units	Percentage
Dynamic?	Yes
Validation	None
When to Change	Increasing this value on systems with multiple CPUs and lots of memory, which are subject to intense periods of memory demand, enables the pager to spend more time attempting to find memory.
Commitment Level	Unstable

## **handspreadpages**

Description	The Solaris OS uses a two-handed clock algorithm to look for pages that are candidates for reclaiming when memory is low. The first hand of the clock walks through memory marking pages as unused. The second hand walks through memory some distance after the first hand,
-------------	--

	checking to see if the page is still marked as unused. If so, the page is subject to being reclaimed. The distance between the first hand and the second hand is <code>handsreadpages</code> .
Data Type	Unsigned long
Default	<code>fastscan</code>
Range	1 to maximum number of physical memory pages on the system
Units	Pages
Dynamic?	Yes. This parameter requires that the kernel <code>reset_hands</code> parameter also be set to a non-zero value. Once the new value of <code>handsreadpages</code> has been recognized, <code>reset_hands</code> is set to zero.
Validation	The value is set to the lesser of either the amount of physical memory and the <code>handsreadpages</code> <i>value</i> .
When to Change	When you want to increase the amount of time that pages are potentially resident before being reclaimed. Increasing this value increases the separation between the hands, and therefore, the amount of time before a page can be reclaimed.
Commitment Level	Unstable

## **pages\_before\_pager**

Description	Defines part of a system threshold that immediately frees pages after an I/O completes instead of storing the pages for possible reuse. The threshold is <code>lotsfree + pages_before_pager</code> . The NFS environment also uses this threshold to curtail its asynchronous activities as memory pressure mounts.
Data Type	Signed integer
Default	200
Range	1 to amount of physical memory
Units	Pages
Dynamic?	No
Validation	None
When to Change	You might change this parameter when the majority of I/O is done for pages that are truly read or written once and never referenced again. Setting this variable to a larger amount of memory keeps adding pages to the free list.

You might also change this parameter when the system is subject to bursts of severe memory pressure. A larger value here helps maintain a larger cushion against the pressure.

Commitment Level Unstable

## maxpgio

**Description** Defines the maximum number of page I/O requests that can be queued by the paging system. This number is divided by 4 to get the actual maximum number used by the paging system. This parameter is used to throttle the number of requests as well as to control process swapping.

**Data Type** Signed integer

**Default** 40

**Range** 1 to a variable maximum that depends on the system architecture, but mainly by the I/O subsystem, such as the number of controllers, disks, and disk swap size

**Units** I/Os

**Dynamic?** No

**Validation** None

**Implicit** The maximum number of I/O requests from the pager is limited by the size of a list of request buffers, which is currently sized at 256.

**When to Change** Increase this parameter to page out memory faster. A larger value might help to recover faster from memory pressure if more than one swap device is configured or if the swap device is a striped device. Note that the existing I/O subsystem should be able to handle the additional I/O load. Also, increased swap I/O could degrade application I/O performance if the swap partition and application files are on the same disk.

Commitment Level Unstable

## Swapping-Related Parameters

Swapping in the Solaris OS is accomplished by the swapfs pseudo file system. The combination of space on swap devices and physical memory is treated as the pool of space available to support the system for maintaining backing store for anonymous memory. The system attempts to allocate space from disk devices first, and then uses physical memory as backing store. When swapfs is forced to use system memory for backing store, limits are enforced to ensure that the system does not deadlock because of excessive consumption by swapfs.

### swapfs\_reserve

Description	Defines the amount of system memory that is reserved for use by system (UID = 0) processes.
Data Type	Unsigned long
Default	The smaller of 4 MB and 1/16th of physical memory
Range	The minimum value is 4 MB or 1/16th of physical memory, whichever is smaller, expressed as pages using the page size returned by <code>getpagesize</code> .  The maximum value is the number of physical memory pages. The maximum value should be no more than 10 percent of physical memory. The system does not enforce this range, other than that described in the Validation section.
Units	Pages
Dynamic?	No
Validation	None
When to Change	Generally not necessary. Only change when recommended by a software provider, or when system processes are terminating because of an inability to obtain swap space. A much better solution is to add physical memory or additional swap devices to the system.
Commitment Level	Unstable

### swapfs\_minfree

Description	Defines the desired amount of physical memory to be kept free for the rest of the system. Attempts to reserve memory for use as swap space by any process that causes the system's perception of available memory to
-------------	--

	fall below this value are rejected. Pages reserved in this manner can only be used for locked-down allocations by the kernel or by user-level processes.
Data Type	Unsigned long
Default	The larger of 2 MB and 1/8th of physical memory
Range	1 to amount of physical memory
Units	Pages
Dynamic?	No
Validation	None
When to Change	When processes are failing because of an inability to obtain swap space, yet the system has memory available.
Commitment Level	Unstable

## Kernel Memory Allocator

The Solaris kernel memory allocator distributes chunks of memory for use by clients inside the kernel. The allocator creates a number of caches of varying size for use by its clients. Clients can also request the allocator to create a cache for use by that client (for example, to allocate structures of a particular size). Statistics about each cache that the allocator manages can be seen by using the `ksstat -c kmem_cache` command.

Occasionally, systems might panic because of memory corruption. The kernel memory allocator supports a debugging interface (a set of flags), that performs various integrity checks on the buffers. The kernel memory allocator also collects information on the allocators. The integrity checks provide the opportunity to detect errors closer to where they actually occurred. The collected information provides additional data for support people when they try to ascertain the reason for the panic.

Use of the flags incurs additional overhead and memory usage during system operations. The flags should only be used when a memory corruption problem is suspected.

### **kmem\_flags**

Description	The Solaris kernel memory allocator has various debugging and test options that were extensively used during the internal development cycle of the Solaris OS. Starting in the Solaris 2.5 release, a subset of these options became available. They are controlled by the <code>kmem_flags</code> variable, which was set with a kernel debugger, and then rebooting the
-------------	---

system. Because of issues with the timing of the instantiation of the kernel memory allocator and the parsing of the `/etc/system` file, it was not possible to set these flags in the `/etc/system` file until the Solaris 8 release.

Five supported flag settings are described here.

Flag	Setting	Description
AUDIT	0x1	The allocator maintains a log that contains recent history of its activity. The number of items logged depends on whether CONTENTS is also set. The log is a fixed size. When space is exhausted, earlier records are reclaimed.
TEST	0x2	The allocator writes a pattern into freed memory and checks that the pattern is unchanged when the buffer is next allocated. If some portion of the buffer is changed, then the memory was probably used by a client that had previously allocated and freed the buffer. If an overwrite is identified, the system panics.
REDZONE	0x4	The allocator provides extra memory at the end of the requested buffer and inserts a special pattern into that memory. When the buffer is freed, the pattern is checked to see if data was written past the end of the buffer. If an overwrite is identified, the kernel panics.
CONTENTS	0x8	The allocator logs up to 256 bytes of buffer contents when the buffer is freed. This flag requires that AUDIT also be set.
		The numeric value of these flags can be logically added together and set by the <code>/etc/system</code> file, starting in the Solaris 8 release, or for previous releases, by booting <code>kadb</code> and setting the flags before starting the kernel.
LITE	0x100	Does minimal integrity checking when a buffer is allocated and freed. When enabled, the allocator checks that the redzone has not been written into, that a freed buffer is not being freed again, and that the buffer being freed is the size that was allocated. This flag is available as of the Solaris 7 3/99 release. Do not combine this flag with any other flags.



Data Type	Signed integer
Default	0 (disabled)
Range	0 (disabled) or 1 - 15 or 256 (0x100)
Dynamic?	Yes. Changes made during runtime only affect new kernel memory caches. After system initialization, the creation of new caches is rare.
Validation	None
When to Change	When memory corruption is suspected
Commitment Level	Unstable

## **kmem\_stackinfo**

Description	<p>If the <code>kmem_stackinfo</code> variable is enabled in the <code>/etc/system</code> file at kernel thread creation time, the kernel thread stack is filled with a specific pattern instead of filled with zeros. During kernel thread execution, this kernel thread stack pattern is progressively overwritten. A simple count from the stack top until the pattern is not found gives a high watermark value, which is the maximum kernel stack space used by a kernel thread. This mechanism allows the following features:</p> <ul style="list-style-type: none"> <li>▪ Compute the percentage of kernel thread stack really used (a high watermark) for current kernel threads in the system</li> <li>▪ When a kernel thread ends, the system logs the last kernel threads that have used the most of their kernel thread stacks before dying to a small circular memory buffer</li> </ul>
Data Type	Unsigned integer
Default	0 (disabled)
Range	0 (disabled) or 1 (enabled)
Dynamic?	Yes
Validation	None
When to Change	When you want to monitor kernel thread stack usage. Keep in mind that when <code>kmem_stackinfo</code> is enabled, the performance of creating and deleting <code>kthreads</code> is decreased. For more information, see the <a href="#">Chapter 5, “Built-In Commands,” in <i>Oracle Solaris Modular Debugger Guide</i></a> .
Zone Configuration	This parameter must be set in the global zone.
Commitment Level	Unstable

# General Driver Parameters

## moddebug

Description	When this parameter is enabled, messages about various steps in the module loading process are displayed.
Data Type	Signed integer
Default	0 (messages off)
Range	<p>Here are the most useful values:</p> <ul style="list-style-type: none"> <li> <b>0x80000000</b> – Prints [un] loading... message. For every module loaded, messages such as the following appear on the console and in the /var/adm/messages file:           <pre>Nov 5 16:12:28 sys genunix: [ID 943528 kern.notice] load 'sched/TS_DPTBL' id 9 loaded @ 0x10126438/ 0x10438dd8 size 132/2064 Nov 5 16:12:28 sys genunix: [ID 131579 kern.notice] installing TS_DPTBL, module id 9.</pre> </li> <li> <b>0x40000000</b> – Prints detailed error messages. For every module loaded, messages such as the following appear on the console and in the /var/adm/messages file:           <pre>Nov 5 16:16:50 sys krtld: [ID 284770 kern.notice] kobj_open: can't open /platform/SUNW,Ultra-80/kernel/ sched/TS_DPTBL Nov 5 16:16:50 sys krtld: [ID 284770 kern.notice] kobj_open: can't open /platform/sun4u/kernel/sched/ TS_DPTBL Nov 5 16:16:50 sys krtld: [ID 797908 kern.notice] kobj_open: '/kernel/sch... Nov 5 16:16:50 sys krtld: [ID 605504 kern.notice] descr = 0x2a Nov 5 16:16:50 sys krtld: [ID 642728 kern.notice] kobj_read_file: size=34, Nov 5 16:16:50 sys krtld: [ID 217760 kern.notice] offset=0 Nov 5 16:16:50 sys krtld: [ID 136382 kern.notice] kobj_read: req 8192 bytes, Nov 5 16:16:50 sys krtld: [ID 295989 kern.notice] got 4224 Nov 5 16:16:50 sys krtld: [ID 426732 kern.notice] read 1080 bytes Nov 5 16:16:50 sys krtld: [ID 720464 kern.notice] copying 34 bytes Nov 5 16:16:50 sys krtld: [ID 234587 kern.notice] count = 34 [33 lines elided] Nov 5 16:16:50 sys genunix: [ID 943528 kern.notice]</pre> </li> </ul>

```
load 'sched/TS_DPTBL' id 9 loaded @ 0x10126438/
0x10438dd8 size 132/2064
Nov 5 16:16:50 sys genunix: [ID 131579 kern.notice]
installing TS_DPTBL, module id 9.
Nov 5 16:16:50 sys genunix: [ID 324367 kern.notice]
init 'sched/TS_DPTBL' id 9 loaded @ 0x10126438/
0x10438dd8 size 132/2064
```

- 0x20000000 - Prints even more detailed messages. This value doesn't print any additional information beyond what the 0x40000000 flag does during system boot. However, this value does print additional information about releasing the module when the module is unloaded.

These values can be added together to set the final value.

Dynamic?	Yes
Validation	None
When to Change	When a module is either not loading as expected, or the system seems to hang while loading modules. Note that when 0x40000000 is set, system boot is slowed down considerably by the number of messages written to the console.
Commitment Level	Unstable

## ddi\_msix\_alloc\_limit

Description	This parameter, available on x86 systems only, controls the number of Extended Message Signaled Interrupts (MSI-X) that a device instance can allocate. Due to an existing system limitation, the default value is 2. You can increase the number of MSI-X interrupts that a device instance can allocate by increasing the value of this parameter. This parameter can be set either by editing the <code>/etc/system</code> file or by setting it with <code>mdb</code> before the device driver attach occurs.
Data Type	Signed integer
Default	2
Range	1 to 16
Dynamic?	Yes
Validation	None
When to Change	To increase the number of MSI-X interrupts that a device instance can allocate. However, if you increase the number of MSI-X interrupts that a device instance can allocate, adequate interrupts might not be

available to satisfy all allocation requests. If this happens, some devices might stop functioning or the system might fail to boot. Reduce the value or remove the parameter in this case.

Commitment Level    Unstable

## Network Driver Parameters

### **igb Parameters**

#### **mr\_enable**

Description	This parameter enables or disables multiple receive and transmit queues that are used by the <code>igb</code> network driver. This parameter can be set by editing the <code>/kernel/drv/igb.conf</code> file before the <code>igb</code> driver attach occurs.
Data Type	Boolean
Default	1 (disable multiple queues)
Range	0 (enable multiple queues) or 1 (disable multiple queues)
Dynamic?	No
Validation	None
When to Change	To enable or disable multiple receive and transmit queues that are used by the <code>igb</code> network driver.
Commitment Level	Unstable

#### **intr\_force**

Description	This parameter is used to force an interrupt type, such as MSI, MSI-X, or legacy, that is used by the <code>igb</code> network driver. This parameter can be set by editing the <code>/kernel/drv/igb.conf</code> file before the <code>igb</code> driver attach occurs.
Data Type	Unsigned integer
Default	0 (do not force an interrupt type)
Range	0 (do not force an interrupt type) 1 (force MSI-X interrupt type)

	2 (force MSI interrupt type)
	3 (force legacy interrupt type)
Dynamic?	No
Validation	None
When to Change	To force an interrupt type that is used by the igb network driver.
Commitment Level	Unstable

## ixgbe Parameters

### tx\_queue\_number

Description	This parameter controls the number of transmit queues that are used by the ixgbe network driver. You can increase the number of transmit queues by increasing the value of this parameter. This parameter can be set by editing the <code>/kernel/drv/ixgbe.conf</code> file before the ixgbe driver attach occurs.
Data Type	Unsigned integer
Default	8
Range	1 to 32
Dynamic?	No
Validation	None
When to Change	To change the number of transmit queues that are used by the ixgbe network driver.
Commitment Level	Unstable

### rx\_queue\_number

Description	This parameter controls the number of receive queues that are used by the ixgbe network driver. You can increase the number of receive queues by increasing the value of this parameter. This parameter can be set by editing the <code>/kernel/drv/ixgbe.conf</code> file before the ixgbe driver attach occurs.
Data Type	Unsigned integer
Default	8

Range	1 to 64
Dynamic?	No
Validation	None
When to Change	To change the number of receive queues that are used by the ixgbe network driver.
Commitment Level	Unstable

### **intr\_throttling**

Description	This parameter controls the interrupt throttling rate of the ixgbe network driver. You can increase the rate of interrupt by decreasing the value of this parameter. This parameter can be set by editing the <code>/kernel/drv/ixgbe.conf</code> file before the ixgbe driver attach occurs.
Data Type	Unsigned integer
Default	200
Range	0 to 65535
Dynamic?	No
Validation	None
When to Change	To change the interrupt throttling rate that is used by the ixgbe network driver.
Commitment Level	Unstable

### **rx\_limit\_per\_intr**

Description	This parameter controls the maximum number of receive queue buffer descriptors per interrupt that are used by the ixgbe network driver. You can increase the number of receive queue buffer descriptors by increasing the value of this parameter. This parameter can be set by editing the <code>/kernel/drv/ixgbe.conf</code> file before the ixgbe driver attach occurs.
Data Type	Unsigned integer
Default	256
Range	16 to 4096
Dynamic?	No
Validation	None

When to Change To change the number of receive queue buffer descriptors that are handled per interrupt by the ixgbe network driver.

Commitment Level Unstable

### **tx\_ring\_size**

Description This parameter controls the transmit queue size that is used by the ixgbe network driver. You can increase the transmit queue size by increasing the value of this parameter. This parameter can be set by editing the `/kernel/drv/ixgbe.conf` file before the ixgbe driver attach occurs.

Data Type Unsigned integer

Default 1024

Range 64 to 4096

Dynamic? No

Validation None

When to Change To change the transmit queue size that is used by the ixgbe network driver.

Commitment Level Unstable

### **rx\_ring\_size**

Description This parameter controls the receive queue size that is used by the ixgbe network driver. You can increase the receive queue size by increasing the value of this parameter. This parameter can be set by editing the `/kernel/drv/ixgbe.conf` file before the ixgbe driver attach occurs.

Data Type Unsigned integer

Default 1024

Range 64 to 4096

Dynamic? No

Validation None

When to Change To change the receive queue size that is used by the ixgbe network driver.

Commitment Level Unstable

**tx\_copy\_threshold**

Description	This parameter controls the transmit buffer copy threshold that is used by the ixgbe network driver. You can increase the transmit buffer copy threshold by increasing the value of this parameter. This parameter can be set by editing the <code>/kernel/drv/ixgbe.conf</code> file before the ixgbe driver attach occurs.
Data Type	Unsigned integer
Default	512
Range	0 to 9126
Dynamic?	No
Validation	None
When to Change	To change the transmit buffer copy threshold that is used by the ixgbe network driver.
Commitment Level	Unstable

**rx\_copy\_threshold**

Description	This parameter controls the receive buffer copy threshold that is used by the ixgbe network driver. You can increase the receive buffer copy threshold by increasing the value of this parameter. This parameter can be set by editing the <code>/kernel/drv/ixgbe.conf</code> file before the ixgbe driver attach occurs.
Data Type	Unsigned integer
Default	128
Range	0 to 9126
Dynamic?	No
Validation	None
When to Change	To change the receive buffer copy threshold that is used by the ixgbe network driver.
Commitment Level	Unstable



## General I/O Parameters

### maxphys

Description	Defines the maximum size of physical I/O requests. If a driver encounters a request larger than this size, the driver breaks the request into maxphys sized chunks. File systems can and do impose their own limit.
Data Type	Signed integer
Default	131,072 (sun4u or sun4v) or 57,344 (x86). The sd driver uses the value of 1,048,576 if the drive supports wide transfers. The ssd driver uses 1,048,576 by default.
Range	Machine-specific page size to MAXINT
Units	Bytes
Dynamic?	Yes, but many file systems load this value into a per-mount point data structure when the file system is mounted. A number of drivers load the value at the time a device is attached to a driver-specific data structure.
Validation	None
When to Change	<p>When doing I/O to and from raw devices in large chunks. Note that a DBMS doing OLTP operations issues large numbers of small I/Os. Changing maxphys does not result in any performance improvement in that case.</p> <p>You might also consider changing this parameter when doing I/O to and from a UFS file system where large amounts of data (greater than 64 KB) are being read or written at any one time. The file system should be optimized to increase contiguity. For example, increase the size of the cylinder groups and decrease the number of inodes per cylinder group. UFS imposes an internal limit of 1 MB on the maximum I/O size it transfers.</p>
Commitment Level	Unstable

### rlim\_fd\_max

Description	Specifies the “hard” limit on file descriptors that a single process might have open. Overriding this limit requires superuser privilege.
-------------	---

Data Type	Signed integer
Default	65,536
Range	1 to MAXINT
Units	File descriptors
Dynamic?	No
Validation	None
When to Change	<p>When the maximum number of open files for a process is not enough. Other limitations in system facilities can mean that a larger number of file descriptors is not as useful as it might be. For example:</p> <ul style="list-style-type: none"><li>▪ A 32-bit program using standard I/O is limited to 256 file descriptors. A 64-bit program using standard I/O can use up to 2 billion descriptors. Specifically, standard I/O refers to the <code>stdio(3C)</code> functions in <code>libc(3LIB)</code>.</li><li>▪ <code>select</code> is by default limited to 1024 descriptors per <code>fd_set</code>. For more information, see <code>select(3C)</code>. Starting with the Solaris 7 release, 32-bit application code can be recompiled with a larger <code>fd_set</code> size (less than or equal to 65,536). A 64-bit application uses an <code>fd_set</code> size of 65,536, which cannot be changed.</li></ul> <p>An alternative to changing this on a system wide basis is to use the <code>plimit(1)</code> command. If a parent process has its limits changed by <code>plimit</code>, all children inherit the increased limit. This alternative is useful for daemons such as <code>inetd</code>.</p>
Commitment Level	Unstable

## `rlim_fd_cur`

Description	Defines the “soft” limit on file descriptors that a single process can have open. A process might adjust its file descriptor limit to any value up to the “hard” limit defined by <code>rlim_fd_max</code> by using the <code>setrlimit()</code> call or by issuing the <code>limit</code> command in whatever shell it is running. You do not require superuser privilege to adjust the limit to any value less than or equal to the hard limit.
Data Type	Signed integer
Default	256
Range	1 to MAXINT

Units	File descriptors
Dynamic?	No
Validation	Compared to <code>rlim_fd_max</code> . If <code>rlim_fd_cur</code> is greater than <code>rlim_fd_max</code> , <code>rlim_fd_cur</code> is reset to <code>rlim_fd_max</code> .
When to Change	When the default number of open files for a process is not enough. Increasing this value means only that it might not be necessary for a program to use <code>setrlimit</code> to increase the maximum number of file descriptors available to it.
Commitment Level	Unstable

## General File System Parameters

### **ncsize**

Description	<p>Defines the number of entries in the directory name look-up cache (DNLC). This parameter is used by UFS, NFS, and ZFS to cache elements of path names that have been resolved.</p> <p>Starting with the Solaris 8 6/00 release, the DNLC also caches negative look-up information, which means it caches a name not found in the cache.</p>
Data Type	Signed integer
Default	$(4 \times (v.v\_proc + \text{maxusers}) + 320) + (4 \times (v.v\_proc + \text{maxusers}) + 320) / 100$
Range	0 to MAXINT
Units	DNLC entries
Dynamic?	No
Validation	None. Larger values cause the time it takes to unmount a file system to increase as the cache must be flushed of entries for that file system during the unmount process.
When to Change	Prior to the Solaris 8 6/00 release, it was difficult to determine whether the cache was too small. You could make this inference by noting the number of entries returned by <code>kstat -n ncstats</code> . If the number seems high, given the system workload and file access pattern, this might be due to the size of the DNLC.

Starting with the Solaris 8 6/00 release, you can use the `kstat -n dnlcstats` command to determine when entries have been removed from the DNLC because it was too small. The sum of the `pick_heuristic` and the `pick_last` parameters represents otherwise valid entries that were reclaimed because the cache was too small.

Excessive values of `ncsize` have an immediate impact on the system because the system allocates a set of data structures for the DNLC based on the value of `ncsize`. A system running a 32-bit kernel allocates 36-byte structures for `ncsize`, while a system running a 64-bit kernel allocates 64-byte structures for `ncsize`. The value has a further effect on UFS and NFS, unless `ufs_ninode` and `nfs:nrnode` are explicitly set.

Commitment Level Unstable

## **rstchown**

Description Indicates whether the POSIX semantics for the `chown` system call are in effect. POSIX semantics are as follows:

- A process cannot change the owner of a file, unless it is running with UID 0.
- A process cannot change the group ownership of a file to a group in which it is not currently a member, unless it is running as UID 0.

For more information, see [chown\(2\)](#).

Data Type Signed integer

Default 1, indicating that POSIX semantics are used

Range 0 = POSIX semantics not in force or 1 = POSIX semantics used

Units Toggle (on/off)

Dynamic? Yes

Validation None

When to Change When POSIX semantics are not wanted. Note that turning off POSIX semantics opens the potential for various security holes. Doing so also opens the possibility of a user changing ownership of a file to another user and being unable to retrieve the file without intervention from the user or the system administrator.

Commitment Level Obsolete

## dnlc\_dir\_enable

Description Enables large directory caching

---

**Note** – This parameter has no effect on NFS or ZFS file systems.

---

Data Type Unsigned integer

Default 1 (enabled)

Range 0 (disabled) or 1 (enabled)

Dynamic? Yes, but do not change this tunable dynamically. You can enable this parameter if it was originally disabled. Or, you can disable this parameter if it was originally enabled. However, enabling, disabling, and then enabling this parameter might lead to stale directory caches.

Validation No

When to Change Directory caching has no known problems. However, if problems occur, then set `dnlc_dir_enable` to 0 to disable caching.

Commitment Level Unstable

## dnlc\_dir\_min\_size

Description Specifies the minimum number of entries cached for one directory.

---

**Note** – This parameter has no effect on NFS or ZFS file systems.

---

Data Type Unsigned integer

Default 40

Range 0 to MAXUINT (no maximum)

Units Entries

Dynamic? Yes, this parameter can be changed at any time.

Validation None

When to Change If performance problems occur with caching small directories, then increase `dnlc_dir_min_size`. Note that individual file systems might have their own range limits for caching directories. For instance, UFS

limits directories to a minimum of `ufs_min_dir_cache` bytes (approximately 1024 entries), assuming 16 bytes per entry.

Commitment Level    Unstable

## **dnlc\_dir\_max\_size**

Description                Specifies the maximum number of entries cached for one directory.

---

**Note** – This parameter has no effect on NFS or ZFS file systems.

---

Data Type                 Unsigned integer

Default                    MAXUINT (no maximum)

Range                      0 to MAXUINT

Dynamic?                 Yes, this parameter can be changed at any time.

Validation                None

When to Change            If performance problems occur with large directories, then decrease `dnlc_dir_max_size`.

Commitment Level        Unstable

## **segmap\_percent**

Description                Defines the maximum amount of memory that is used for the fast-access file system cache. This pool of memory is subtracted from the free memory list.

Data Type                 Unsigned integer

Default                    12 percent of free memory at system startup time

Range                      2 MB to 100 percent of `physmem`

Units                      % of physical memory

Dynamic?                 No

Validation                None

When to Change            If heavy file system activity is expected, and sufficient free memory is available, you should increase the value of this parameter.

Commitment Level Unstable

## UFS Parameters

### bufhwm and bufhwm\_pct

Description	<p>Defines the maximum amount of memory for caching I/O buffers. The buffers are used for writing file system metadata (superblocks, inodes, indirect blocks, and directories). Buffers are allocated as needed until the amount of memory (in KB) to be allocated exceed <code>bufhwm</code>. At this point, metadata is purged from the buffer cache until enough buffers are reclaimed to satisfy the request.</p> <p>For historical reasons, <code>bufhwm</code> does not require the <code>ufs:</code> prefix.</p>
Data Type	Signed integer
Default	2 percent of physical memory
Range	80 KB to 20 percent of physical memory, or 2 TB, whichever is less. Consequently, <code>bufhwm_pct</code> can be between 1 and 20.
Units	<p><code>bufhwm</code>: KB</p> <p><code>bufhwm_pct</code>: percent of physical memory</p>
Dynamic?	<p>No. <code>bufhwm</code> and <code>bufhwm_pct</code> are only evaluated at system initialization to compute hash bucket sizes. The limit in bytes calculated from these parameters is then stored in a data structure that adjusts this value as buffers are allocated and deallocated.</p> <p>Attempting to adjust this value without following the locking protocol on a running system can lead to incorrect operation.</p> <p>Modifying <code>bufhwm</code> or <code>bufhwm_pct</code> at runtime has no effect.</p>
Validation	<p>If <code>bufhwm</code> is less than its lower limit of 80 KB or greater than its upper limit (the lesser of 20 percent of physical memory, 2 TB, or one quarter (1/4) of the maximum amount of kernel heap), it is reset to the upper limit. The following message appears on the system console and in the <code>/var/adm/messages</code> file if an invalid value is attempted:</p> <pre>"binit: bufhwm (value attempted) out of range (range start..range end). Using N as default."</pre>

“Value attempted” refers to the value specified in the `/etc/system` file or by using a kernel debugger. *N* is the value computed by the system based on available system memory.

Likewise, if `bufhwm_pct` is set to a value that is outside the allowed range of 1 percent to 20 percent, it is reset to the default of 2 percent. And, the following message appears on the system console and in the `/var/adm/messages` file:

```
"binit: bufhwm_pct(value attempted) out of range(0..20).
        Using 2 as default."
```

If both `bufhwm` or `bufhwm_pct` are set to non-zero values, `bufhwm` takes precedence.

#### When to Change

Because buffers are only allocated as they are needed, the overhead from the default setting is the required allocation of control structures for the buffer hash headers. These structures consume 52 bytes per potential buffer on a 32-bit kernel and 96 bytes per potential buffer on a 64-bit kernel.

On a 512-MB 64-bit kernel, the number of hash chains calculates to  $10316 / 32 == 322$ , which scales up to next power of 2, 512. Therefore, the hash headers consume  $512 \times 96$  bytes, or 48 KB. The hash header allocations assume that buffers are 32 KB.

The amount of memory, which has not been allocated in the buffer pool, can be found by looking at the `bfreelist` structure in the kernel with a kernel debugger. The field of interest in the structure is `b_bufsize`, which is the possible remaining memory in bytes. Looking at it with the `buf` macro by using the `mdb` command:

```
# mdb -k
Loading modules: [ unix krtld genunix ip nfs ipc ]
> bfreelist::print "struct buf" b_bufsize
b_bufsize = 0x225800
```

The default value for `bufhwm` on this system, with 6 GB of memory, is 122277. You cannot determine the number of header structures used because the actual buffer size requested is usually larger than 1 KB. However, some space might be profitably reclaimed from control structure allocation for this system.

The same structure on a 512-MB system shows that only 4 KB of 10144 KB has not been allocated. When the `biostats kstat` is examined with `kstat -n biostats`, it is determined that the system had a reasonable ratio of `buffer_cache_hits` to `buffer_cache_lookups` as well. As such, the default setting is reasonable for that system.



Commitment Level      Unstable

## ndquot

Description	Defines the number of quota structures for the UFS file system that should be allocated. Relevant only if quotas are enabled on one or more UFS file systems. Because of historical reasons, the <code>ufs:</code> prefix is not needed.
Data Type	Signed integer
Default	$((\text{maxusers} \times 40) / 4) + \text{max\_nprocs}$
Range	0 to MAXINT
Units	Quota structures
Dynamic?	No
Validation	None. Excessively large values hang the system.
When to Change	When the default number of quota structures is not enough. This situation is indicated by the following message displayed on the console or written in the message log:  <code>dquot table full</code>
Commitment Level	Unstable

## ufs\_ninode

Description	<p>Specifies the number of inodes to be held in memory. Inodes are cached globally for UFS, not on a per-file system basis.</p> <p>A key parameter in this situation is <code>ufs_ninode</code>. This parameter is used to compute two key limits that affect the handling of inode caching. A high watermark of <math>\text{ufs\_ninode} / 2</math> and a low watermark of <math>\text{ufs\_ninode} / 4</math> are computed.</p> <p>When the system is done with an inode, one of two things can happen:</p> <ul style="list-style-type: none"> <li>▪ The file referred to by the inode is no longer on the system so the inode is deleted. After it is deleted, the space goes back into the inode cache for use by another inode (which is read from disk or created for a new file).</li> </ul>
-------------	---

- The file still exists but is no longer referenced by a running process. The inode is then placed on the idle queue. Any referenced pages are still in memory.

When inodes are idled, the kernel defers the idling process to a later time. If a file system is a logging file system, the kernel also defers deletion of inodes. Two kernel threads handle this deferred processing. Each thread is responsible for one of the queues.

When the deferred processing is done, the system drops the inode onto either a delete queue or an idle queue, each of which has a thread that can run to process it. When the inode is placed on the queue, the queue occupancy is checked against the low watermark. If the queue occupancy exceeds the low watermark, the thread associated with the queue is awakened. After the queue is awakened, the thread runs through the queue and forces any pages associated with the inode out to disk and frees the inode. The thread stops when it has removed 50 percent of the inodes on the queue at the time it was awakened.

A second mechanism is in place if the idle thread is unable to keep up with the load. When the system needs to find a vnode, it goes through the `ufs_vget` routine. The *first* thing `vget` does is check the length of the idle queue. If the length is above the high watermark, then it takes two inodes off the idle queue and “idles” them (flushes pages and frees inodes). `vget` does this *before* it gets an inode for its own use.

The system does attempt to optimize by placing inodes with no in-core pages at the head of the idle list and inodes with pages at the end of the idle list. However, the system does no other ordering of the list. Inodes are always removed from the front of the idle queue.

The only time that inodes are removed from the queues as a whole is when a synchronization, unmount, or remount occur.

For historical reasons, this parameter does not require the `ufs :` prefix.

Data Type	Signed integer
Default	<code>ncsize</code>
Range	0 to <code>MAXINT</code>
Units	Inodes
Dynamic?	Yes
Validation	If <code>ufs_ninode</code> is less than or equal to zero, the value is set to <code>ncsize</code> .

When to Change	When the default number of inodes is not enough. If the <code>maxsize</code> reached field as reported by <code>kstat -n inode_cache</code> is larger than the <code>maxsize</code> field in the <code>kstat</code> , the value of <code>ufs_ninode</code> might be too small. Excessive inode idling can also be a problem.  You can identify excessive inode idling by using <code>kstat -n inode_cache</code> to look at the <code>inode_cache</code> <code>kstat</code> . Thread idles are inodes idled by the background threads while <code>vget idles</code> are idles by the requesting process before using an inode.
Commitment Level	Unstable

## **ufs\_WRITES**

Description	If <code>ufs_WRITES</code> is non-zero, the number of bytes outstanding for writes on a file is checked. See <code>ufs_HW</code> to determine whether the write should be issued or deferred until only <code>ufs_LW</code> bytes are outstanding. The total number of bytes outstanding is tracked on a per-file basis so that if the limit is passed for one file, it won't affect writes to other files.
Data Type	Signed integer
Default	1 (enabled)
Range	0 (disabled) or 1 (enabled)
Units	Toggle (on/off)
Dynamic?	Yes
Validation	None
When to Change	When you want UFS write throttling turned off entirely. If sufficient I/O capacity does not exist, disabling this parameter can result in long service queues for disks.
Commitment Level	Unstable

## **ufs\_LW and ufs\_HW**

Description	<code>ufs_HW</code> specifies the number of bytes outstanding on a single file barrier value. If the number of bytes outstanding is greater than this value and <code>ufs_WRITES</code> is set, then the write is deferred. The write is deferred by putting the thread issuing the write to sleep on a condition variable.
-------------	---

`ufs_LW` is the barrier for the number of bytes outstanding on a single file below which the condition variable on which other sleeping processes are toggled. When a write completes and the number of bytes is less than `ufs_LW`, then the condition variable is toggled, which causes all threads waiting on the variable to awaken and try to issue their writes.

Data Type	Signed integer
Default	8 x 1024 x 1024 for <code>ufs_LW</code> and 16 x 1024 x 1024 for <code>ufs_HW</code>
Range	0 to MAXINT
Units	Bytes
Dynamic?	Yes
Validation	None
Implicit	<code>ufs_LW</code> and <code>ufs_HW</code> have meaning only if <code>ufs_WRITES</code> is not equal to zero. <code>ufs_HW</code> and <code>ufs_LW</code> should be changed together to avoid needless churning when processes awaken and find that either they cannot issue a write (when <code>ufs_LW</code> and <code>ufs_HW</code> are too close) or they might have waited longer than necessary (when <code>ufs_LW</code> and <code>ufs_HW</code> are too far apart).
When to Change	Consider changing these values when file systems consist of striped volumes. The aggregate bandwidth available can easily exceed the current value of <code>ufs_HW</code> . Unfortunately, this parameter is not a per-file system setting.  You might also consider changing this parameter when <code>ufs_throttles</code> is a non-trivial number. Currently, <code>ufs_throttles</code> can only be accessed with a kernel debugger.
Commitment Level	Unstable

## freebehind

Description	Enables the <code>freebehind</code> algorithm. When this algorithm is enabled, the system bypasses the file system cache on newly read blocks when sequential I/O is detected during times of heavy memory use.
Data Type	Boolean
Default	1 (enabled)
Range	0 (disabled) or 1 (enabled)

Dynamic?	Yes
Validation	None
When to Change	The <code>freebehind</code> algorithm can occur too easily. If no significant sequential file system activity is expected, disabling <code>freebehind</code> makes sure that all files, no matter how large, will be candidates for retention in the file system page cache. For more fine-grained tuning, see <code>smallfile</code> .
Commitment Level	Unstable

## **smallfile**

Description	<p>Determines the size threshold of files larger than this value are candidates for no cache retention under the <code>freebehind</code> algorithm.</p> <p>Large memory systems contain enough memory to cache thousands of 10-MB files without making severe memory demands. However, this situation is highly application dependent.</p> <p>The goal of the <code>smallfile</code> and <code>freebehind</code> parameters is to reuse cached information, without causing memory shortfalls by caching too much.</p>
Data Type	Signed integer
Default	32,768
Range	0 to 2,147,483,647
Dynamic?	Yes
Validation	None
When to Change	Increase <code>smallfile</code> if an application does sequential reads on medium-sized files and can most likely benefit from buffering, and the system is not otherwise under pressure for free memory. Medium-sized files are 32 KB to 2 GB in size.
Commitment Level	Unstable

## TMPFS Parameters

### **tmpfs:tmpfs\_maxkmem**

Description	Defines the maximum amount of kernel memory that TMPFS can use for its data structures (tmpnodes and directory entries).
Data Type	Unsigned long
Default	One page or 4 percent of physical memory, whichever is greater.
Range	Number of bytes in one page (8192 for sun4u or sun4v systems, 4096 for all other systems) to 25 percent of the available kernel memory at the time TMPFS was first used.
Units	Bytes
Dynamic?	Yes
Validation	None
When to Change	Increase if the following message is displayed on the console or written in the messages file:  <code>tmp_memalloc: tmpfs over memory limit</code>  The current amount of memory used by TMPFS for its data structures is held in the <code>tmp_kmemspace</code> field. This field can be examined with a kernel debugger.
Commitment Level	Unstable

### **tmpfs:tmpfs\_minfree**

Description	Defines the minimum amount of swap space that TMPFS leaves for the rest of the system.
Data Type	Signed long
Default	256
Range	0 to maximum swap space size
Units	Pages
Dynamic?	Yes
Validation	None

When to Change	To maintain a reasonable amount of swap space on systems with large amounts of TMPFS usage, you can increase this number. The limit has been reached when the console or messages file displays the following message:  <i>fs-name: File system full, swap space limit exceeded</i>
Commitment Level	Unstable

## Pseudo Terminals

Pseudo terminals, ptys, are used for two purposes in Solaris software:

- Supporting remote logins by using the `telnet`, `rlogin`, or `rsh` commands
- Providing the interface through which the X Window system creates command interpreter windows

The default number of pseudo-terminals is sufficient for a desktop workstation. So, tuning focuses on the number of ptys available for remote logins.

Previous versions of Solaris required that steps be taken to explicitly configure the system for the preferred number of ptys. Starting with the Solaris 8 release, a new mechanism removes the necessity for tuning in most cases. The default number of ptys is now based on the amount of memory on the system. This default should be changed only to restrict or increase the number of users who can log in to the system.

Three related variables are used in the configuration process:

- `pt_cnt` – Default maximum number of ptys.
- `pt_pctofmem` – Percentage of kernel memory that can be dedicated to pty support structures. A value of zero means that no remote users can log in to the system.
- `pt_max_pty` – Hard maximum for number of ptys.

`pt_cnt` has a default value of zero, which tells the system to limit logins based on the amount of memory specified in `pt_pctofmem`, unless `pt_max_pty` is set. If `pt_cnt` is non-zero, ptys are allocated until this limit is reached. When that threshold is crossed, the system looks at `pt_max_pty`. If `pt_max_pty` has a non-zero value, it is compared to `pt_cnt`. The pty allocation is allowed if `pt_cnt` is less than `pt_max_pty`. If `pt_max_pty` is zero, `pt_cnt` is compared to the number of ptys supported based on `pt_pctofmem`. If `pt_cnt` is less than this value, the pty allocation is allowed. Note that the limit based on `pt_pctofmem` only comes into play if both `pt_cnt` and `ptms_ptymax` have default values of zero.

To put a hard limit on ptys that is different than the maximum derived from `pt_pctofmem`, set `pt_cnt` and `ptms_ptymax` in `/etc/system` to the preferred number of ptys. The setting of `ptms_pctofmem` is not relevant in this case.

To dedicate a different percentage of system memory to pty support and let the operating system manage the explicit limits, do the following:

- Do not set `pt_cnt` or `ptms_ptymax` in `/etc/system`.
- Set `pt_pctofmem` in `/etc/system` to the preferred percentage. For example, set `pt_pctofmem=10` for a 10 percent setting.

Note that the memory is not actually allocated until it is used in support of a pty. Once memory is allocated, it remains allocated.

## pt\_cnt

Description	The number of available <code>/dev/pts</code> entries is dynamic up to a limit determined by the amount of physical memory available on the system. <code>pt_cnt</code> is one of three variables that determines the minimum number of logins that the system can accommodate. The default maximum number of <code>/dev/pts</code> devices the system can support is determined at boot time by computing the number of pty structures that can fit in a percentage of system memory (see <code>pt_pctofmem</code> ). If <code>pt_cnt</code> is zero, the system allocates up to that maximum. If <code>pt_cnt</code> is non-zero, the system allocates to the greater of <code>pt_cnt</code> and the default maximum.
Data Type	Unsigned integer
Default	0
Range	0 to <code>maxpid</code>
Units	Logins/windows
Dynamic?	No
Validation	None
When to Change	When you want to explicitly control the number of users who can remotely log in to the system.
Commitment Level	Unstable

## pt\_pctofmem

Description	Specifies the maximum percentage of physical memory that can be consumed by data structures to support <code>/dev/pts</code> entries. A system running a 64-bit kernel consumes 176 bytes per <code>/dev/pts</code> entry. A system running a 32-bit kernel consumes 112 bytes per <code>/dev/pts</code> entry.
-------------	---



Data Type	Unsigned integer
Default	5
Range	0 to 100
Units	Percentage
Dynamic?	No
Validation	None
When to Change	When you want to either restrict or increase the number of users who can log in to the system. A value of zero means that no remote users can log in to the system.
Commitment Level	Unstable

## **pt\_max\_pty**

Description	Defines the maximum number of ptys the system offers
Data Type	Unsigned integer
Default	0 (Uses system-defined maximum)
Range	0 to MAXUINT
Units	Logins/windows
Dynamic?	Yes
Validation	None
Implicit	Should be greater than or equal to pt_cnt. Value is not checked until the number of ptys allocated exceeds the value of pt_cnt.
When to Change	When you want to place an absolute ceiling on the number of logins supported, even if the system could handle more based on its current configuration values.
Commitment Level	Unstable

## STREAMS Parameters

### nstrpush

Description	Specifies the number of modules that can be inserted into (pushed onto) a STREAM.
Data Type	Signed integer
Default	9
Range	9 to 16
Units	Modules
Dynamic?	Yes
Validation	None
When to Change	At the direction of your software vendor. No messages are displayed when a STREAM exceeds its permitted push count. A value of EINVAL is returned to the program that attempted the push.
Commitment Level	Unstable

### strmsgsz

Description	Specifies the maximum number of bytes that a single system call can pass to a STREAM to be placed in the data part of a message. Any <code>write</code> exceeding this size is broken into multiple messages. For more information, see <a href="#">write(2)</a> .
Data Type	Signed integer
Default	65,536
Range	0 to 262,144
Units	Bytes
Dynamic?	Yes
Validation	None
When to Change	When <code>putmsg</code> calls return ERANGE. For more information, see <a href="#">putmsg(2)</a> .
Commitment Level	Unstable

## strctlsz

Description	Specifies the maximum number of bytes that a single system call can pass to a STREAM to be placed in the control part of a message
Data Type	Signed integer
Default	1024
Range	0 to MAXINT
Units	Bytes
Dynamic?	Yes
Validation	None
When to Change	At the direction of your software vendor. <code>putmsg(2)</code> calls return <code>ERANGE</code> if they attempt to exceed this limit.
Commitment Level	Unstable

## System V Message Queues

System V message queues provide a message-passing interface that enables the exchange of messages by queues created in the kernel. Interfaces are provided in the Solaris environment to enqueue and dequeue messages. Messages can have a type associated with them. Enqueueing places messages at the end of a queue. Dequeueing removes the first message of a specific type from the queue or the first message if no type is specified.

For detailed information on tuning these system resources, see [Chapter 6, “Resource Controls \(Overview\)”](#), in *System Administration Guide: Oracle Solaris Zones, Oracle Solaris 10 Containers, and Resource Management*.

## System V Semaphores

System V semaphores provide counting semaphores in the Solaris OS. A *semaphore* is a counter used to provide access to a shared data object for multiple processes. In addition to the standard set and release operations for semaphores, System V semaphores can have values that are incremented and decremented as needed (for example, to represent the number of resources available). System V semaphores also provide the ability to do operations on a group of semaphores simultaneously as well as to have the system undo the last operation by a process if the process dies.

## System V Shared Memory

System V shared memory allows the creation of a segment by a process. Cooperating processes can attach to the memory segment (subject to access permissions on the segment) and gain access to the data contained in the segment. This capability is implemented as a loadable module. Entries in the `/etc/system` file must contain the `shmsys:` prefix. Starting with the Solaris 7 release, the `keyserv` daemon uses System V shared memory.

A special kind of shared memory known as *intimate shared memory* (ISM) is used by DBMS vendors to maximize performance. When a shared memory segment is made into an ISM segment, the memory for the segment is locked. This feature enables a faster I/O path to be followed and improves memory usage. A number of kernel resources describing the segment are then shared between all processes that attach to the segment in ISM mode.

### `segspt_minfree`

Description	Identifies pages of system memory that cannot be allocated for ISM shared memory.
Data Type	Unsigned long
Default	5 percent of available system memory when the first ISM segment is created
Range	0 to 50 percent of physical memory
Units	Pages
Dynamic?	Yes
Validation	None. Values that are too small can cause the system to hang or performance to severely degrade when memory is consumed with ISM segments.
When to Change	On database servers with large amounts of physical memory using ISM, the value of this parameter can be decreased. If ISM segments are not used, this parameter has no effect. A maximum value of 128 MB (0x4000) is almost certainly sufficient on large memory machines.
Commitment Level	Unstable

## Scheduling

### rechoose\_interval

Description	Specifies the number of clock ticks before a process is deemed to have lost all affinity for the last CPU it ran on. After this interval expires, any CPU is considered a candidate for scheduling a thread. This parameter is relevant only for threads in the timesharing class. Real-time threads are scheduled on the first available CPU.
Data Type	Signed integer
Default	3
Range	0 to MAXINT
Dynamic?	Yes
Validation	None
When to Change	When caches are large, or when the system is running a critical process or a set of processes that seem to suffer from excessive cache misses not caused by data access patterns.  Consider using the processor set capabilities available as of the Solaris 2.6 release or processor binding before changing this parameter. For more information, see <a href="#">psrset(1M)</a> or <a href="#">pbind(1M)</a> .
Commitment Level	Unstable

## Timers

### hires\_tick

Description	When set, this parameter causes the Solaris OS to use a system clock rate of 1000 instead of the default value of 100.
Data Type	Signed integer
Default	0
Range	0 (disabled) or 1 (enabled)
Dynamic?	No. Causes new system timing variable to be set at boot time. Not referenced after boot.

Validation	None
When to Change	When you want timeouts with a resolution of less than 10 milliseconds, and greater than or equal to 1 millisecond.
Commitment Level	Unstable

## **timer\_max**

Description	Specifies the number of POSIX timers available.
Data Type	Signed integer
Default	32
Range	0 to MAXINT
Dynamic?	No. Increasing the value can cause a system crash.
Validation	None
When to Change	When the default number of timers offered by the system is inadequate. Applications receive an EAGAIN error when executing <code>timer_create</code> system calls.
Commitment Level	Unstable

## **sun4u or sun4v Specific Parameters**

### **consistent\_coloring**

Description	<p>Starting with the Solaris 2.6 release, the ability to use different page placement policies on the UltraSPARC (sun4u) platform was introduced. A page placement policy attempts to allocate physical page addresses to maximize the use of the L2 cache. Whatever algorithm is chosen as the default algorithm, that algorithm can potentially provide less optimal results than another algorithm for a particular application set. This parameter changes the placement algorithm selected for all processes on the system.</p> <p>Based on the size of the L2 cache, memory is divided into bins. The page placement code allocates a page from a bin when a page fault first occurs on an unmapped page. The page chosen depends on which of the three possible algorithms are used:</p>
-------------	---

- Page coloring – Various bits of the virtual address are used to determine the bin from which the page is selected. This is the default algorithm in the Solaris 8 release. `consistent_coloring` is set to zero to use this algorithm. No per-process history exists for this algorithm.
- Virtual addr=physical address – Consecutive pages in the program selects pages from consecutive bins. `consistent_coloring` is set to 1 to use this algorithm. No per-process history exists for this algorithm.
- Bin-hopping – Consecutive pages in the program generally allocate pages from every other bin, but the algorithm occasionally skips more bins. `consistent_coloring` is set to 2 to use this algorithm. Each process starts at a randomly selected bin, and a per-process memory of the last bin allocated is kept.

Dynamic?	Yes
Validation	None. Values larger than 2 cause a number of <code>WARNING: AS_2_BIN: bad consistent coloring value</code> messages to appear on the console. The system hangs immediately thereafter. A power-cycle is required to recover.
When to Change	When the primary workload of the system is a set of long-running high-performance computing (HPC) applications. Changing this value might provide better performance. File servers, database servers, and systems with a number of active processes (for example, compile or time sharing servers) do not benefit from changes.
Commitment Level	Unstable

## **tsb\_alloc\_hiwater\_factor**

Description      Initializes `tsb_alloc_hiwater` to impose an upper limit on the amount of physical memory that can be allocated for translation storage buffers (TSBs) as follows:

$$\text{tsb\_alloc\_hiwater} = \text{physical memory (bytes)} / \text{tsb\_alloc\_hiwater\_factor}$$

When the memory that is allocated to TSBs is equal to the value of `tsb_alloc_hiwater`, the TSB memory allocation algorithm attempts to reclaim TSB memory as pages are unmapped.

Exercise caution when using this factor to increase the value of `tsb_alloc_hiwater`. To prevent system hangs, the resulting high water value must be considerably lower than the value of `swapfs_minfree` and `segspt_minfree`.

Data Type	Integer
Default	32
Range	1 to MAXINIT

Note that a factor of 1 makes all physical memory available for allocation to TSBs, which could cause the system to hang. A factor that is too high will not leave memory available for allocation to TSBs, decreasing system performance.

Dynamic?	Yes
Validation	None
When to Change	Change the value of this parameter if the system has many processes that attach to very large shared memory segments. Under most circumstances, tuning of this variable is not necessary.
Commitment Level	Unstable

## **default\_tsb\_size**

Description	Selects size of the initial translation storage buffers (TSBs) allocated to all processes.
Data Type	Integer
Default	Default is 0 (8 KB), which corresponds to 512 entries
Range	Possible values are:

Value	Description
0	8 KB
1	16 KB
3	32 KB
4	128 KB
5	256 KB
6	512 KB



Value	Description
7	1 Mbyte

Dynamic?	Yes
Validation	None
When to Change	Generally, you do not need to change this value. However, doing so might provide some advantages if the majority of processes on the system have a larger than average working set, or if resident set size (RSS) sizing is disabled.
Commitment Level	Unstable

## **enable\_tsb\_rss\_sizing**

Description	Enables a resident set size (RSS) based TSB sizing heuristic.
Data Type	Boolean
Default	1 (TSBs can be resized)
Range	0 (TSBs remain at <code>tsb_default_size</code> ) or 1 (TSBs can be resized)
	If set to 0, then <code>tsb_rss_factor</code> is ignored.
Dynamic?	Yes
Validation	Yes
When to Change	Can be set to 0 to prevent growth of the TSBs. Under most circumstances, this parameter should be left at the default setting.
Commitment Level	Unstable

## **tsb\_rss\_factor**

Description	Controls the RSS to TSB span ratio of the RSS sizing heuristic. This factor divided by 512 yields the percentage of the TSB span which must be resident in memory before the TSB is considered as a candidate for resizing.
Data Type	Integer

Default	384, resulting in a value of 75%. Thus, when the TSB is 3/4 full, its size will be increased. Note that some virtual addresses typically map to the same slot in the TSB. Therefore, conflicts can occur before the TSB is at 100% full.
Range	0 to 512
Dynamic?	Yes
Validation	None
When to Change	<p>If the system is experiencing an excessive number of traps due to TSB misses, for example, due to virtual address conflicts in the TSB, you might consider decreasing this value toward 0.</p> <p>For example, changing <code>tsb_rss_factor</code> to 256 (effectively, 50%) instead of 384 (effectively, 75%) might help eliminate virtual address conflicts in the TSB in some cases, but will use more kernel memory, particularly on a heavily loaded system.</p> <p>TSB activity can be monitored with the <code>trapstat -T</code> command.</p>
Commitment Level	Unstable

## Locality Group Parameters

This section provides generic memory tunables, which apply to any SPARC or x86 system that uses a Non-Uniform Memory Architecture (NUMA).

### `lpg_alloc_prefer`

Description	<p>Controls a heuristic for allocation of large memory pages when the requested page size is not immediately available in the local memory group, but could be satisfied from a remote memory group.</p> <p>By default, the Solaris OS allocates a remote large page if local free memory is fragmented, but remote free memory is not. Setting this parameter to 1 indicates that additional effort should be spent attempting to allocate larger memory pages locally, potentially moving smaller pages around to coalesce larger pages in the local memory group.</p>
Data Type	Boolean

Default	0 (Prefer remote allocation if local free memory is fragmented and remote free memory is not)
Range	0 (Prefer remote allocation if local free memory is fragmented and remote free memory is not)  1 (Prefer local allocation whenever possible, even if local free memory is fragmented and remote free memory is not)
Dynamic?	No
Validation	None
When to Change	<p>This parameter might be set to 1 if long-running programs on the system tend to allocate memory that is accessed by a single program, or if memory that is accessed by a group of programs is known to be running in the same locality group (lgroup). In these circumstances, the extra cost of page coalesce operations can be amortized over the long run of the programs.</p> <p>This parameter might be left at the default value (0) if multiple programs tend to share memory across different locality groups, or if pages tend to be used for short periods of time. In these circumstances, quick allocation of the requested size tends to be more important than allocation in a particular location.</p> <p>Page locations and sizes might be observed by using the NUMA observability tools, available at <a href="http://hub.opensolaris.org/bin/view/Main/">http://hub.opensolaris.org/bin/view/Main/</a>. TLB miss activity might be observed by using the <code>trapstat -T</code> command.</p>
Commitment Level	Uncommitted

## **lgrp\_mem\_default\_policy**

Description	This variable reflects the default memory allocation policy used by the Solaris OS. This variable is an integer, and its value should correspond to one of the policies listed in the <code>sys/lgrp.h</code> file.
Data Type	Integer
Default	1, <code>LGRP_MEM_POLICY_NEXT</code> indicating that memory allocation defaults to the home lgroup of the thread performing the memory allocation.
Range	Possible values are:

Value	Description	Comment
0	LGRP_MEM_POLICY_DEFAULT	use system default policy
1	LGRP_MEM_POLICY_NEXT	next to allocating thread's home lgroup
2	LGRP_MEM_POLICY_RANDOM_PROC	randomly across process
3	LGRP_MEM_POLICY_RANDOM_PSET	randomly across processor set
4	LGRP_MEM_POLICY_RANDOM	randomly across all lgroups
5	LGRP_MEM_POLICY_ROUNDROBIN	round robin across all lgroups
6	LGRP_MEM_POLICY_NEXT_CPU	near next CPU to touch memory

Dynamic?	No
Validation	None
When to Change	For applications that are sensitive to memory latencies due to allocations that occur from remote versus local memory on systems that use NUMA.
Commitment Level	Uncommitted

## **lgrp\_mem\_pset\_aware**

Description	<p>If a process is running within a user processor set, this variable determines whether <i>randomly</i> placed memory for the process is selected from among all the lgroups in the system or only from those lgroups that are spanned by the processors in the processor set.</p> <p>For more information about creating processor sets, see <a href="#">psrset(1M)</a>.</p>
Data Type	Boolean
Default	0, the Solaris OS selects memory from all the lgroups in the system
Range	<ul style="list-style-type: none"> <li>▪ 0, the Solaris OS selects memory from all the lgroups in the system (default)</li> <li>▪ 1, try selecting memory only from those lgroups that are spanned by the processors in the processor set. If the first attempt fails, memory can be allocated in any lgroup.</li> </ul>
Dynamic?	No

Validation	None
When to Change	Setting this value to a value of one (1) might lead to more reproducible performance when processor sets are used to isolate applications from one another.
Commitment Level	Uncommitted



# NFS Tunable Parameters

---

This section describes the NFS tunable parameters.

- [“Tuning the NFS Environment” on page 95](#)
- [“NFS Module Parameters” on page 96](#)
- [“nfsrv Module Parameters” on page 123](#)
- [“rpcmod Module Parameters” on page 126](#)

## Where to Find Tunable Parameter Information

Tunable Parameter	For Information
Solaris kernel tunables	<a href="#">Chapter 2, “Oracle Solaris Kernel Tunable Parameters”</a>
Internet Protocol Suite tunable parameters	<a href="#">Chapter 4, “Internet Protocol Suite Tunable Parameters”</a>
Network Cache and Accelerator (NCA) tunable parameters	<a href="#">Chapter 5, “Network Cache and Accelerator Tunable Parameters”</a>

## Tuning the NFS Environment

You can define NFS parameters in the `/etc/system` file, which is read during the boot process. Each parameter includes the name of its associated kernel module. For more information, see [“Tuning a Solaris System” on page 19](#).



**Caution** – The names of the parameters, the modules that they reside in, and the default values can change between releases. Check the documentation for the version of the active SunOS release before making changes or applying values from previous releases.

---

## NFS Module Parameters

This section describes parameters related to the NFS kernel module.

### **nfs:nfs3\_pathconf\_disable\_cache**

Description	Controls the caching of pathconf information for NFS Version 3 mounted file systems.
Data Type	Integer (32-bit)
Default	0 (caching enabled)
Range	0 (caching enabled) or 1 (caching disabled)
Units	Boolean values
Dynamic?	Yes
Validation	None
When to Change	The pathconf information is cached on a per file basis. However, if the server can change the information for a specific file dynamically, use this parameter to disable caching. There is no mechanism for the client to validate its cache entry.
Commitment Level	Unstable

### **nfs:nfs4\_pathconf\_disable\_cache**

Description	Controls the caching of pathconf information for NFS Version 4 mounted file systems.
Data Type	Integer (32-bit)
Default	0 (caching enabled)
Range	0 (caching enabled) or 1 (caching disabled)
Units	Boolean values
Dynamic?	Yes



Validation	None
When to Change	The <code>pathconf</code> information is cached on a per file basis. However, if the server can change the information for a specific file dynamically, use this parameter to disable caching. There is no mechanism for the client to validate its cache entry.
Commitment Level	Unstable

## **nfs:nfs\_allow\_preepoch\_time**

**Description** Controls whether files with incorrect or *negative* time stamps should be made visible on the client.

Historically, neither the NFS client nor the NFS server would do any range checking on the file times being returned. The over-the-wire timestamp values are unsigned and 32-bits long. So, all values have been legal.

However, on a system running a 32-bit Solaris kernel, the timestamp values are signed and 32-bits long. Thus, it would be possible to have a timestamp representation that appeared to be prior to January 1, 1970, or *pre-epoch*.

The problem on a system running a 64-bit Solaris kernel is slightly different. The timestamp values on the 64-bit Solaris kernel are signed and 64-bits long. It is impossible to determine whether a time field represents a full 32-bit time or a negative time, that is, a time prior to January 1, 1970.

It is impossible to determine whether to sign extend a time value when converting from 32 bits to 64 bits. The time value should be sign extended if the time value is truly a negative number. However, the time value should not be sign extended if it does truly represent a full 32-bit time value. This problem is resolved by simply disallowing full 32-bit time values.

Data Type	Integer (32-bit)
Default	0 (32-bit time stamps disabled)
Range	0 (32-bit time stamps disabled) or 1 (32-bit time stamps enabled)
Units	Boolean values
Dynamic?	Yes

Validation	None
When to Change	Even during normal operation, it is possible for the timestamp values on some files to be set very far in the future or very far in the past. If access to these files is preferred using NFS mounted file systems, set this parameter to 1 to allow the timestamp values to be passed through unchecked.
Commitment Level	Unstable

## **nfs:nfs\_cots\_timeo**

Description	Controls the default RPC timeout for NFS version 2 mounted file systems using connection-oriented transports such as TCP for the transport protocol.
Data Type	Signed integer (32-bit)
Default	600 (60 seconds)
Range	0 to $2^{31} - 1$
Units	10th of seconds
Dynamic?	Yes, but the RPC timeout for a file system is set when the file system is mounted. To affect a particular file system, unmount and mount the file system after changing this parameter.
Validation	None
When to Change	TCP does a good job ensuring requests and responses are delivered appropriately. However, if the round-trip times are very large in a particularly slow network, the NFS version 2 client might time out prematurely.  Increase this parameter to prevent the client from timing out incorrectly. The range of values is very large, so increasing this value too much might result in situations where a retransmission is not detected for long periods of time.
Commitment Level	Unstable

## **nfs:nfs3\_cots\_timeo**

Description	Controls the default RPC timeout for NFS version 3 mounted file systems using connection-oriented transports such as TCP for the transport protocol.
-------------	--

Data Type	Signed integer (32-bit)
Default	600 (60 seconds)
Range	0 to $2^{31} - 1$
Units	10th of seconds
Dynamic?	Yes, but the RPC timeout for a file system is set when the file system is mounted. To affect a particular file system, unmount and mount the file system after changing this parameter.
Validation	None
When to Change	TCP does a good job ensuring requests and responses are delivered appropriately. However, if the round-trip times are very large in a particularly slow network, the NFS version 3 client might time out prematurely.  Increase this parameter to prevent the client from timing out incorrectly. The range of values is very large, so increasing this value too much might result in situations where a retransmission is not detected for long periods of time.
Commitment Level	Unstable

## **nfs:nfs4\_cots\_timeo**

Description	Controls the default RPC timeout for NFS version 4 mounted file systems using connection-oriented transports such as TCP for the transport protocol.  The NFS Version 4 protocol specification disallows retransmission over the same TCP connection. Thus, this parameter primarily controls how quickly the client responds to certain events, such as detecting a forced unmount operation or detecting how quickly the server fails over to a new server.
Data Type	Signed integer (32-bit)
Default	600 (60 seconds)
Range	0 to $2^{31} - 1$
Units	10th of seconds
Dynamic?	Yes, but this parameter is set when the file system is mounted. To affect a particular file system, unmount and mount the file system after changing this parameter.

Validation	None
When to Change	<p>TCP does a good job ensuring requests and responses are delivered appropriately. However, if the round-trip times are very large in a particularly slow network, the NFS version 4 client might time out prematurely.</p> <p>Increase this parameter to prevent the client from timing out incorrectly. The range of values is very large, so increasing this value too much might result in situations where a retransmission is not detected for long periods of time.</p>
Commitment Level	Unstable

## **nfs:nfs\_do\_symlink\_cache**

Description	Controls whether the contents of symbolic link files are cached for NFS version 2 mounted file systems.
Data Type	Integer (32-bit)
Default	1 (caching enabled)
Range	0 (caching disabled) or 1 (caching enabled)
Units	Boolean values
Dynamic?	Yes
Validation	None
When to Change	<p>If a server changes the contents of a symbolic link file without updating the modification timestamp on the file or if the granularity of the timestamp is too large, then changes to the contents of the symbolic link file might not be visible on the client for extended periods. In this case, use this parameter to disable the caching of symbolic link contents. Doing so makes the changes immediately visible to applications running on the client.</p>
Commitment Level	Unstable

## **nfs:nfs3\_do\_symlink\_cache**

Description	Controls whether the contents of symbolic link files are cached for NFS version 3 mounted file systems.
Data Type	Integer (32-bit)

Default	1 (caching enabled)
Range	0 (caching disabled) or 1 (caching enabled)
Units	Boolean values
Dynamic?	Yes
Validation	None
When to Change	If a server changes the contents of a symbolic link file without updating the modification timestamp on the file or if the granularity of the timestamp is too large, then changes to the contents of the symbolic link file might not be visible on the client for extended periods. In this case, use this parameter to disable the caching of symbolic link contents. Doing so makes the changes immediately visible to applications running on the client.
Commitment Level	Unstable

## **nfs:nfs4\_do\_symlink\_cache**

Description	Controls whether the contents of symbolic link files are cached for NFS version 4 mounted file systems.
Data Type	Integer (32-bit)
Default	1 (caching enabled)
Range	0 (caching disabled) or 1 (caching enabled)
Units	Boolean values
Dynamic?	Yes
Validation	None
When to Change	If a server changes the contents of a symbolic link file without updating the modification timestamp on the file or if the granularity of the timestamp is too large, then changes to the contents of the symbolic link file might not be visible on the client for extended periods. In this case, use this parameter to disable the caching of symbolic link contents. Doing so makes the changes immediately visible to applications running on the client.
Commitment Level	Unstable

## **nfs:nfs\_dynamic**

Description	Controls whether a feature known as <i>dynamic retransmission</i> is enabled for NFS version 2 mounted file systems using connectionless transports such as UDP. This feature attempts to reduce retransmissions by monitoring server response times and then adjusting RPC timeouts and read- and write- transfer sizes.
Data Type	Integer (32-bit)
Default	1 (enabled)
Range	0 (disabled) or 1 (enabled)
Dynamic?	Yes, but this parameter is set per file system at mount time. To affect a particular file system, unmount and mount the file system after changing this parameter.
Validation	None
When to Change	Do not change this parameter.
Commitment Level	Unstable

## **nfs:nfs3\_dynamic**

Description	Controls whether a feature known as <i>dynamic retransmission</i> is enabled for NFS version 3 mounted file systems using connectionless transports such as UDP. This feature attempts to reduce retransmissions by monitoring server response times and then adjusting RPC timeouts and read- and write- transfer sizes.
Data Type	Integer (32-bit)
Default	0 (disabled)
Range	0 (disabled) or 1 (enabled)
Units	Boolean values
Dynamic?	Yes, but this parameter is set per file system at mount time. To affect a particular file system, unmount and mount the file system after changing this parameter.
Validation	None
When to Change	Do not change this parameter.
Commitment Level	Unstable

## nfs:nfs\_lookup\_neg\_cache

Description	Controls whether a negative name cache is used for NFS version 2 mounted file systems. This negative name cache records file names that were looked up, but not found. The cache is used to avoid over-the-network look-up requests made for file names that are already known to not exist.
Data Type	Integer (32-bit)
Default	1 (enabled)
Range	0 (disabled) or 1 (enabled)
Units	Boolean values
Dynamic?	Yes
Validation	None
When to Change	<p>For the cache to perform correctly, negative entries must be strictly verified before they are used. This consistency mechanism is relaxed slightly for read-only mounted file systems. It is assumed that the file system on the server is not changing or is changing very slowly, and that it is okay for such changes to propagate slowly to the client. The consistency mechanism becomes the normal attribute cache mechanism in this case.</p> <p>If file systems are mounted read-only on the client, but are expected to change on the server and these changes need to be seen immediately by the client, use this parameter to disable the negative cache.</p> <p>If you disable the <code>nfs:nfs_disable_rmdir_cache</code> parameter, you should probably also disable this parameter. For more information, see <a href="#">“nfs:nfs_disable_rmdir_cache” on page 113</a>.</p>
Commitment Level	Unstable

## nfs:nfs3\_lookup\_neg\_cache

Description	Controls whether a negative name cache is used for NFS version 3 mounted file systems. This negative name cache records file names that were looked up, but were not found. The cache is used to avoid over-the-network look-up requests made for file names that are already known to not exist.
Data Type	Integer (32-bit)

Default	1 (enabled)
Range	0 (disabled) or 1 (enabled)
Units	Boolean values
Dynamic?	Yes
Validation	None
When to Change	<p>For the cache to perform correctly, negative entries must be strictly verified before they are used. This consistency mechanism is relaxed slightly for read-only mounted file systems. It is assumed that the file system on the server is not changing or is changing very slowly, and that it is okay for such changes to propagate slowly to the client. The consistency mechanism becomes the normal attribute cache mechanism in this case.</p> <p>If file systems are mounted read-only on the client, but are expected to change on the server and these changes need to be seen immediately by the client, use this parameter to disable the negative cache.</p> <p>If you disable the <code>nfs:nfs_disable_rmdir_cache</code> parameter, you should probably also disable this parameter. For more information, see <a href="#">“nfs:nfs_disable_rmdir_cache” on page 113</a>.</p>
Commitment Level	Unstable

## **nfs:nfs4\_lookup\_neg\_cache**

Description	Controls whether a negative name cache is used for NFS version 4 mounted file systems. This negative name cache records file names that were looked up, but were not found. The cache is used to avoid over-the-network look-up requests made for file names that are already known to not exist.
Data Type	Integer (32-bit)
Default	1 (enabled)
Range	0 (disabled) or 1 (enabled)
Units	Boolean values
Dynamic?	Yes
Validation	None



When to Change	<p>For the cache to perform correctly, negative entries must be strictly verified before they are used. This consistency mechanism is relaxed slightly for read-only mounted file systems. It is assumed that the file system on the server is not changing or is changing very slowly, and that it is okay for such changes to propagate slowly to the client. The consistency mechanism becomes the normal attribute cache mechanism in this case.</p> <p>If file systems are mounted read-only on the client, but are expected to change on the server and these changes need to be seen immediately by the client, use this parameter to disable the negative cache.</p> <p>If you disable the <code>nfs:nfs_disable_rmdir_cache</code> parameter, you should probably also disable this parameter. For more information, see <a href="#">“nfs:nfs_disable_rmdir_cache” on page 113</a>.</p>
Commitment Level	Unstable

## **nfs:nfs\_max\_threads**

Description	<p>Controls the number of kernel threads that perform asynchronous I/O for the NFS version 2 client. Because NFS is based on RPC and RPC is inherently synchronous, separate execution contexts are required to perform NFS operations that are asynchronous from the calling thread.</p> <p>The operations that can be executed asynchronously are read for read-ahead, readdir for readdir read-ahead, write for putpage and pageio operations, commit, and inactive for cleanup operations that the client performs when it stops using a file.</p>
Data Type	Integer (16-bit)
Default	8
Range	0 to $2^{15} - 1$
Units	Threads
Dynamic?	Yes, but this parameter is set per file system at mount time. To affect a particular file system, unmount and mount the file system after changing this parameter.
Validation	None
When to Change	To increase or reduce the number of simultaneous I/O operations that are outstanding at any given time. For example, for a very low bandwidth network, you might want to decrease this value so that the

NFS client does not overload the network. Alternately, if the network is very high bandwidth, and the client and server have sufficient resources, you might want to increase this value. Doing so can more effectively utilize the available network bandwidth, and the client and server resources.

Commitment Level Unstable

## **nfs:nfs3\_max\_threads**

**Description** Controls the number of kernel threads that perform asynchronous I/O for the NFS version 3 client. Because NFS is based on RPC and RPC is inherently synchronous, separate execution contexts are required to perform NFS operations that are asynchronous from the calling thread.

The operations that can be executed asynchronously are read for read-ahead, readdir for readdir read-ahead, write for putpage and pageio requests, and commit.

**Data Type** Integer (16-bit)

**Default** 8

**Range** 0 to  $2^{15} - 1$

**Units** Threads

**Dynamic?** Yes, but this parameter is set per file system at mount time. To affect a particular file system, unmount and mount the file system after changing this parameter.

**Validation** None

**When to Change** To increase or reduce the number of simultaneous I/O operations that are outstanding at any given time. For example, for a very low bandwidth network, you might want to decrease this value so that the NFS client does not overload the network. Alternately, if the network is very high bandwidth, and the client and server have sufficient resources, you might want to increase this value. Doing so can more effectively utilize the available network bandwidth, and the client and server resources.

Commitment Level Unstable

## **nfs:nfs4\_max\_threads**

Description	Controls the number of kernel threads that perform asynchronous I/O for the NFS version 4 client. Because NFS is based on RPC and RPC is inherently synchronous, separate execution contexts are required to perform NFS operations that are asynchronous from the calling thread.  The operations that can be executed asynchronously are read for read-ahead, write-behind, directory read-ahead, and cleanup operations that the client performs when it stops using a file.
Data Type	Integer (16-bit)
Default	8
Range	0 to $2^{15} - 1$
Units	Threads
Dynamic?	Yes, but this parameter is set per file system at mount time. To affect a particular file system, unmount and mount the file system after changing this parameter.
Validation	None
When to Change	To increase or reduce the number of simultaneous I/O operations that are outstanding at any given time. For example, for a very low bandwidth network, you might want to decrease this value so that the NFS client does not overload the network. Alternately, if the network is very high bandwidth, and the client and server have sufficient resources, you might want to increase this value. Doing so can more effectively utilize the available network bandwidth, and the client and server resources.
Commitment Level	Unstable

## **nfs:nfs\_nra**

Description	Controls the number of read-ahead operations that are queued by the NFS version 2 client when sequential access to a file is discovered. These read-ahead operations increase concurrency and read throughput. Each read-ahead request is generally for one logical block of file data.
Data Type	Integer (32-bit)
Default	4

Range	0 to $2^{31} - 1$
Units	Logical blocks. (See “ <a href="#">nfs:nfs_bsize</a> ” on page 114.)
Dynamic?	Yes
Validation	None
When to Change	To increase or reduce the number of read-ahead requests that are outstanding for a specific file at any given time. For example, for a very low bandwidth network or on a low memory client, you might want to decrease this value so that the NFS client does not overload the network or the system memory. Alternately, if the network is very high bandwidth, and the client and server have sufficient resources, you might want to increase this value. Doing so can more effectively utilize the available network bandwidth, and the client and server resources.
Commitment Level	Unstable

## **nfs:nfs3\_nra**

Description	Controls the number of read-ahead operations that are queued by the NFS version 3 client when sequential access to a file is discovered. These read-ahead operations increase concurrency and read throughput. Each read-ahead request is generally for one logical block of file data.
Data Type	Integer (32-bit)
Default	4
Range	0 to $2^{31} - 1$
Units	Logical blocks. (See “ <a href="#">nfs:nfs3_bsize</a> ” on page 114.)
Dynamic?	Yes
Validation	None
When to Change	To increase or reduce the number of read-ahead requests that are outstanding for a specific file at any given time. For example, for a very low bandwidth network or on a low memory client, you might want to decrease this value so that the NFS client does not overload the network or the system memory. Alternately, if the network is very high bandwidth and the client and server have sufficient resources, you might want to increase this value. Doing so can more effectively utilize the available network bandwidth, and the client and server resources.
Commitment Level	Unstable

## nfs:nfs4\_nra

Description	Controls the number of read-ahead operations that are queued by the NFS version 4 client when sequential access to a file is discovered. These read-ahead operations increase concurrency and read throughput. Each read-ahead request is generally for one logical block of file data.
Data Type	Integer (32-bit)
Default	4
Range	0 to $2^{31} - 1$
Units	Logical blocks. (See “ <a href="#">nfs:nfs4_bsize</a> ” on page 115.)
Dynamic?	Yes
Validation	None
When to Change	To increase or reduce the number of read-ahead requests that are outstanding for a specific file at any given time. For example, for a very low bandwidth network or on a low memory client, you might want to decrease this value so that the NFS client does not overload the network or the system memory. Alternately, if the network is very high bandwidth, and the client and server have sufficient resources, you might want to increase this value. Doing so can more effectively utilize the available network bandwidth, and the client and server resources.
Commitment Level	Unstable

## nfs:nrnode

Description	Controls the size of the rnode cache on the NFS client.  The rnode, used by both NFS version 2, 3, and 4 clients, is the central data structure that describes a file on the NFS client. The rnode contains the file handle that identifies the file on the server. The rnode also contains pointers to various caches used by the NFS client to avoid network calls to the server. Each rnode has a one-to-one association with a vnode. The vnode caches file data.  The NFS client attempts to maintain a minimum number of rnodes to attempt to avoid destroying cached data and metadata. When an rnode is reused or freed, the cached data and metadata must be destroyed.
Data Type	Integer (32-bit)

Default	The default setting of this parameter is 0, which means that the value of <code>nnode</code> should be set to the value of the <code>ncsize</code> parameter. Actually, any non positive value of <code>nnode</code> results in <code>nnode</code> being set to the value of <code>ncsize</code> .
Range	1 to $2^{31} - 1$
Units	rnodes
Dynamic?	No. This value can only be changed by adding or changing the parameter in the <code>/etc/system</code> file, and then rebooting the system.
Validation	The system enforces a maximum value such that the <code>rnode</code> cache can only consume 25 percent of available memory.
When to Change	<p>Because <code>rnodes</code> are created and destroyed dynamically, the system tends to settle upon a <code>nnode</code>-size cache, automatically adjusting the size of the cache as memory pressure on the system increases or as more files are simultaneously accessed. However, in certain situations, you could set the value of <code>nnode</code> if the mix of files being accessed can be predicted in advance. For example, if the NFS client is accessing a few very large files, you could set the value of <code>nnode</code> to a small number so that system memory can cache file data instead of <code>rnodes</code>. Alternately, if the client is accessing many small files, you could increase the value of <code>nnode</code> to optimize for storing file metadata to reduce the number of network calls for metadata.</p> <p>Although it is not recommended, the <code>rnode</code> cache can be effectively disabled by setting the value of <code>nnode</code> to 1. This value instructs the client to only cache 1 <code>rnode</code>, which means that it is reused frequently.</p>
Commitment Level	Unstable

## **nfs:nfs\_shrinkreaddir**

**Description** Some older NFS servers might incorrectly handle NFS version 2 REaddir requests for more than 1024 bytes of directory information. This problem is due to a bug in the server implementation. However, this parameter contains a workaround in the NFS version 2 client.

When this parameter is enabled, the client does not generate a REaddir request for larger than 1024 bytes of directory information. If this parameter is disabled, then the over-the-wire size is set to the lesser of either the size passed in by using the `getdents` system call or by using `NFS_MAXDATA`, which is 8192 bytes. For more information, see [getdents\(2\)](#).

Data Type	Integer (32-bit)
Default	0 (disabled)
Range	0 (disabled) or 1 (enabled)
Units	Boolean values
Dynamic?	Yes
Validation	None
When to Change	Examine the value of this parameter if an older NFS version 2 only server is used and interoperability problems occur when the server tries to read directories. Enabling this parameter might cause a slight decrease in performance for applications that read directories.
Commitment Level	Unstable

## **nfs:nfs3\_shrinkreaddir**

Description	<p>Some older NFS servers might incorrectly handle NFS version 3 READDIR requests for more than 1024 bytes of directory information. This problem is due to a bug in the server implementation. However, this parameter contains a workaround in the NFS version 3 client.</p> <p>When this parameter is enabled, the client does not generate a READDIR request for larger than 1024 bytes of directory information. If this parameter is disabled, then the over-the-wire size is set to the minimum of either the size passed in by using the <code>getdents</code> system call or by using <code>MAXBSIZE</code>, which is 8192 bytes. For more information, see <a href="#">getdents(2)</a>.</p>
Data Type	Integer (32-bit)
Default	0 (disabled)
Range	0 (disabled) or 1 (enabled)
Units	Boolean values
Dynamic?	Yes
Validation	None
When to Change	Examine the value of this parameter if an older NFS version 3 only server is used and interoperability problems occur when the server tries to read directories. Enabling this parameter might cause a slight decrease in performance for applications that read directories.

Commitment Level    Unstable

## **nfs:nfs\_write\_error\_interval**

Description	Controls the time duration in between logging ENOSPC and EDQUOT write errors received by the NFS client. This parameter affects NFS version 2, 3, and 4 clients.
Data Type	Long integer (32 bits on 32-bit platforms and 64 bits on 64-bit platforms)
Default	5 seconds
Range	0 to $2^{31} - 1$ on 32-bit platforms 0 to $2^{63} - 1$ on 64-bit platforms
Units	Seconds
Dynamic?	Yes
Validation	None
When to Change	Increase or decrease the value of this parameter in response to the volume of messages being logged by the client. Typically, you might want to increase the value of this parameter to decrease the number of out of space messages being printed when a full file system on a server is being actively used.
Commitment Level	Unstable

## **nfs:nfs\_write\_error\_to\_cons\_only**

Description	Controls whether NFS write errors are logged to the system console and <code>syslog</code> or to the system console only. This parameter affects messages for NFS version 2, 3, and 4 clients.
Data Type	Integer (32-bit)
Default	0 (system console and <code>syslog</code> )
Range	0 (system console and <code>syslog</code> ) or 1 (system console)
Units	Boolean values
Dynamic?	Yes
Validation	None



**When to Change** Examine the value of this parameter to avoid filling up the file system containing the messages logged by the `syslogd` daemon. When this parameter is enabled, messages are printed on the system console only and are not copied to the `syslog` messages file.

**Commitment Level** Unstable

## **nfs:nfs\_disable\_rddir\_cache**

**Description** Controls the use of a cache to hold responses from `REaddir` and `REaddirplus` requests. This cache avoids over-the-wire calls to the server to retrieve directory information.

**Data Type** Integer (32-bit)

**Default** 0 (caching enabled)

**Range** 0 (caching enabled) or 1 (caching disabled)

**Units** Boolean values

**Dynamic?** Yes

**Validation** None

**When to Change** Examine the value of this parameter if interoperability problems develop due to a server that does not update the modification time on a directory when a file or directory is created in it or removed from it. The symptoms are that new names do not appear in directory listings after they have been added to the directory or that old names do not disappear after they have been removed from the directory.

This parameter controls the caching for NFS version 2, 3, and 4 mounted file systems. This parameter applies to all NFS mounted file systems, so caching cannot be disabled or enabled on a per file system basis.

If you disable this parameter, you should also disable the following parameters to prevent bad entries in the DNLC negative cache:

- “`nfs:nfs_lookup_neg_cache`” on page 103
- “`nfs:nfs3_lookup_neg_cache`” on page 103
- “`nfs:nfs4_lookup_neg_cache`” on page 104

**Commitment Level** Unstable

## **nfs:nfs\_bsize**

Description	Controls the logical block size used by the NFS version 2 client. This block size represents the amount of data that the client attempts to read from or write to the server when it needs to do an I/O.
Data Type	Unsigned integer (32-bit)
Default	8192 bytes
Range	0 to $2^{31} - 1$
Units	Bytes
Dynamic?	Yes, but the block size for a file system is set when the file system is mounted. To affect a particular file system, unmount and mount the file system after changing this parameter.
Validation	None. Setting this parameter too low or too high might cause the system to malfunction. Do not set this parameter to anything less than PAGESIZE for the specific platform. Do not set this parameter too high because it might cause the system to hang while waiting for memory allocations to be granted.
When to Change	Do not change this parameter.
Commitment Level	Unstable

## **nfs:nfs3\_bsize**

Description	Controls the logical block size used by the NFS version 3 client. This block size represents the amount of data that the client attempts to read from or write to the server when it needs to do an I/O.
Data Type	Unsigned integer (32-bit)
Default	32,768 (32 KB)
Range	0 to $2^{31} - 1$
Units	Bytes
Dynamic?	Yes, but the block size for a file system is set when the file system is mounted. To affect a particular file system, unmount and mount the file system after changing this parameter.
Validation	None. Setting this parameter too low or too high might cause the system to malfunction. Do not set this parameter to anything less than

	PAGESIZE for the specific platform. Do not set this parameter too high because it might cause the system to hang while waiting for memory allocations to be granted.
When to Change	Examine the value of this parameter when attempting to change the maximum data transfer size. Change this parameter in conjunction with the <code>nfs:nfs3_max_transfer_size</code> parameter. If larger transfers are preferred, increase both parameters. If smaller transfers are preferred, then just reducing this parameter should suffice.
Commitment Level	Unstable

## **nfs:nfs4\_bsize**

Description	Controls the logical block size used by the NFS version 4 client. This block size represents the amount of data that the client attempts to read from or write to the server when it needs to do an I/O.
Data Type	Unsigned integer (32-bit)
Default	32,768 (32 KB)
Range	0 to $2^{31} - 1$
Units	Bytes
Dynamic?	Yes, but the block size for a file system is set when the file system is mounted. To affect a particular file system, unmount and mount the file system after changing this parameter.
Validation	None. Setting this parameter too low or too high might cause the system to malfunction. Do not set this parameter to anything less than PAGESIZE for the specific platform. Do not set this parameter too high because it might cause the system to hang while waiting for memory allocations to be granted.
When to Change	Examine the value of this parameter when attempting to change the maximum data transfer size. Change this parameter in conjunction with the <code>nfs:nfs4_max_transfer_size</code> parameter. If larger transfers are preferred, increase both parameters. If smaller transfers are preferred, then just reducing this parameter should suffice.
Commitment Level	Unstable

## nfs:nfs\_async\_clusters

Description	<p>Controls the mix of asynchronous requests that are generated by the NFS version 2 client. The four types of asynchronous requests are read-ahead, putpage, pageio, and readdir-ahead. The client attempts to round-robin between these different request types to attempt to be fair and not starve one request type in favor of another.</p> <p>However, the functionality in some NFS version 2 servers such as write gathering depends upon certain behaviors of existing NFS Version 2 clients. In particular, this functionality depends upon the client sending out multiple WRITE requests at about the same time. If one request is taken out of the queue at a time, the client would be defeating this server functionality designed to enhance performance for the client.</p> <p>Thus, use this parameter to control the number of requests of each request type that are sent out before changing types.</p>
Data Type	Unsigned integer (32-bit)
Default	1
Range	0 to $2^{31} - 1$
Units	Asynchronous requests
Dynamic?	Yes, but the cluster setting for a file system is set when the file system is mounted. To affect a particular file system, unmount and mount the file system after changing this parameter.
Validation	None. However, setting the value of this parameter to 0 causes all of the queued requests of a particular request type to be processed before moving on to the next type. This effectively disables the fairness portion of the algorithm.
When to Change	To increase the number of each type of asynchronous request that is generated before switching to the next type. Doing so might help with server functionality that depends upon clusters of requests coming from the client.
Commitment Level	Unstable

## nfs:nfs3\_async\_clusters

Description	Controls the mix of asynchronous requests that are generated by the NFS version 3 client. The five types of asynchronous requests are read-ahead, putpage, pageio, readdir-ahead, and commit. The client
-------------	--

attempts to round-robin between these different request types to attempt to be fair and not starve one request type in favor of another.

However, the functionality in some NFS version 3 servers such as write gathering depends upon certain behaviors of existing NFS version 3 clients. In particular, this functionality depends upon the client sending out multiple WRITE requests at about the same time. If one request is taken out of the queue at a time, the client would be defeating this server functionality designed to enhance performance for the client.

Thus, use this parameter to control the number of requests of each request type that are sent out before changing types.

Data Type	Unsigned integer (32-bit)
Default	1
Range	0 to $2^{31} - 1$
Units	Asynchronous requests
Dynamic?	Yes, but the cluster setting for a file system is set when the file system is mounted. To affect a particular file system, unmount and mount the file system after changing this parameter.
Validation	None. However, setting the value of this parameter to 0 causes all of the queued requests of a particular request type to be processed before moving on to the next type. This value effectively disables the fairness portion of the algorithm.
When to Change	To increase the number of each type of asynchronous operation that is generated before switching to the next type. Doing so might help with server functionality that depends upon clusters of operations coming from the client.
Commitment Level	Unstable

## **nfs:nfs4\_async\_clusters**

**Description** Controls the mix of asynchronous requests that are generated by the NFS version 4 client. The six types of asynchronous requests are read-ahead, putpage, pageio, readdir-ahead, commit, and inactive. The client attempts to round-robin between these different request types to attempt to be fair and not starve one request type in favor of another.

However, the functionality in some NFS version 4 servers such as write gathering depends upon certain behaviors of existing NFS version 4

clients. In particular, this functionality depends upon the client sending out multiple WRITE requests at about the same time. If one request is taken out of the queue at a time, the client would be defeating this server functionality designed to enhance performance for the client.

Thus, use this parameter to control the number of requests of each request type that are sent out before changing types.

Data Type	Unsigned integer (32-bit)
Default	1
Range	0 to $2^{31} - 1$
Units	Asynchronous requests
Dynamic?	Yes, but the cluster setting for a file system is set when the file system is mounted. To affect a particular file system, unmount and mount the file system after changing this parameter.
Validation	None. However, setting the value of this parameter to 0 causes all of the queued requests of a particular request type to be processed before moving on to the next type. This effectively disables the fairness portion of the algorithm.
When to Change	To increase the number of each type of asynchronous request that is generated before switching to the next type. Doing so might help with server functionality that depends upon clusters of requests coming from the client.
Commitment Level	Unstable

## **nfs:nfs\_async\_timeout**

Description	Controls the duration of time that threads, which execute asynchronous I/O requests, sleep with nothing to do before exiting. When there are no more requests to execute, each thread goes to sleep. If no new requests come in before this timer expires, the thread wakes up and exits. If a request does arrive, a thread is woken up to execute requests until there are none again. Then, the thread goes back to sleep waiting for another request to arrive, or for the timer to expire.
Data Type	Integer (32-bit)
Default	6000 (1 minute expressed as 60 sec * 100Hz)
Range	0 to $2^{31} - 1$
Units	Hz. (Typically, the clock runs at 100Hz.)

Dynamic?	Yes
Validation	None. However, setting this parameter to a non positive value causes these threads exit as soon as there are no requests in the queue for them to process.
When to Change	If the behavior of applications in the system is known precisely and the rate of asynchronous I/O requests can be predicted, it might be possible to tune this parameter to optimize performance slightly in either of the following ways: <ul style="list-style-type: none"> <li>▪ By making the threads expire more quickly, thus freeing up kernel resources more quickly</li> <li>▪ By making the threads expire more slowly, thus avoiding thread create and destroy overhead</li> </ul>
Commitment Level	Unstable

## **nfs:nacache**

Description	Tunes the number of hash queues that access the file access cache on the NFS client. The file access cache stores file access rights that users have with respect to files that they are trying to access. The cache itself is dynamically allocated. However, the hash queues used to index into the cache are statically allocated. The algorithm assumes that there is one access cache entry per active file and four of these access cache entries per hash bucket. Thus, by default, the value of this parameter is set to the value of the <code>nnode</code> parameter.
Data Type	Integer (32-bit)
Default	The default setting of this parameter is 0. This value means that the value of <code>nacache</code> should be set to the value of the <code>nnode</code> parameter.
Range	1 to $2^{31} - 1$
Units	Access cache entries
Dynamic?	No. This value can only be changed by adding or changing the parameter in the <code>/etc/system</code> file, and then rebooting system.
Validation	None. However, setting this parameter to a negative value will probably cause the system to try to allocate a very large set of hash queues. While trying to do so, the system is likely to hang.
When to Change	Examine the value of this parameter if the basic assumption of one access cache entry per file would be violated. This violation could occur for systems in a timesharing mode where multiple users are accessing

the same file at about the same time. In this case, it might be helpful to increase the expected size of the access cache so that the hashed access to the cache stays efficient.

Commitment Level Unstable

## **nfs:nfs3\_jukebox\_delay**

Description	Controls the duration of time that the NFS version 3 client waits to transmit a new request after receiving the NFS3ERR_JUKEBOX error from a previous request. The NFS3ERR_JUKEBOX error is generally returned from the server when the file is temporarily unavailable for some reason. This error is generally associated with hierarchical storage, and CD or tape jukeboxes.
Data Type	Long integer (32 bits on 32-bit platforms and 64 bits on 64-bit platforms)
Default	1000 (10 seconds expressed as 10 sec * 100Hz)
Range	0 to $2^{31} - 1$ on 32-bit platforms 0 to $2^{63} - 1$ on 64-bit platforms
Units	Hz. (Typically, the clock runs at 100Hz.)
Dynamic?	Yes
Validation	None
When to Change	Examine the value of this parameter and perhaps adjust it to match the behaviors exhibited by the server. Increase this value if the delays in making the file available are long in order to reduce network overhead due to repeated retransmissions. Decrease this value to reduce the delay in discovering that the file has become available.
Commitment Level	Unstable

## **nfs:nfs3\_max\_transfer\_size**

Description	Controls the maximum size of the data portion of an NFS version 3 READ, WRITE, REaddir, or REaddirplus request. This parameter controls both the maximum size of the request that the server returns as well as the maximum size of the request that the client generates.
Data Type	Integer (32-bit)



Default	1,048,576 (1 Mbyte)
Range	0 to $2^{31} - 1$
Units	Bytes
Dynamic?	Yes, but this parameter is set per file system at mount time. To affect a particular file system, unmount and mount the file system after changing this parameter.
Validation	<p>None. However, setting the maximum transfer size on the server to 0 is likely to cause clients to malfunction or just decide not to attempt to talk to the server.</p> <p>There is also a limit on the maximum transfer size when using NFS over the UDP transport. UDP has a hard limit of 64 KB per datagram. This 64 KB must include the RPC header as well as other NFS information, in addition to the data portion of the request. Setting the limit too high might result in errors from UDP and communication problems between the client and the server.</p>
When to Change	<p>To tune the size of data transmitted over the network. In general, the <code>nfs:nfs3_bsize</code> parameter should also be updated to reflect changes in this parameter.</p> <p>For example, when you attempt to increase the transfer size beyond 32 KB, update <code>nfs:nfs3_bsize</code> to reflect the increased value. Otherwise, no change in the over-the-wire request size is observed. For more information, see <a href="#">“nfs:nfs3_bsize” on page 114</a>.</p> <p>If you want to use a smaller transfer size than the default transfer size, use the mount command's <code>-wsize</code> or <code>-rsize</code> option on a per-file system basis.</p>
Commitment Level	Unstable

## **nfs:nfs4\_max\_transfer\_size**

Description	Controls the maximum size of the data portion of an NFS version 4 READ, WRITE, REaddir, or REaddirplus request. This parameter controls both the maximum size of the request that the server returns as well as the maximum size of the request that the client generates.
Data Type	Integer (32-bit)
Default	32,768 (32 KB)
Range	0 to $2^{31} - 1$

Units	Bytes
Dynamic?	Yes, but this parameter is set per file system at mount time. To affect a particular file system, unmount and mount the file system after changing this parameter.
Validation	None. However, setting the maximum transfer size on the server to 0 is likely to cause clients to malfunction or just decide not to attempt to talk to the server.  There is also a limit on the maximum transfer size when using NFS over the UDP transport. For more information on the maximum for UDP, see “ <a href="#">nfs:nfs3_max_transfer_size</a> ” on page 120.
When to Change	To tune the size of data transmitted over the network. In general, the <code>nfs:nfs4_bsize</code> parameter should also be updated to reflect changes in this parameter.  For example, when you attempt to increase the transfer size beyond 32 KB, update <code>nfs:nfs4_bsize</code> to reflect the increased value. Otherwise, no change in the over-the-wire request size is observed. For more information, see “ <a href="#">nfs:nfs4_bsize</a> ” on page 115.  If you want to use a smaller transfer size than the default transfer size, use the mount command's <code>-wsize</code> or <code>-rsize</code> option on a per-file system basis.
Commitment Level	Unstable

## **nfs:nfs3\_max\_transfer\_size\_clts**

Description	Controls the maximum size of the data portion of an NFS version 3 READ, WRITE, REaddir, or REaddirPLUS request over UDP. This parameter controls both the maximum size of the request that the server returns as well as the maximum size of the request that the client generates.
Data Type	Integer (32-bit)
Default	32,768 (32 KB)
Range	0 to $2^{31} - 1$
Units	Bytes
Dynamic?	Yes, but this parameter is set per file system at mount time. To affect a particular file system, unmount and mount the file system after changing this parameter.

Validation	None. However, setting the maximum transfer size on the server to 0 is likely to cause clients to malfunction or just decide not to attempt to talk to the server.
When to Change	Do not change this parameter.
Commitment Level	Unstable

## **nfs:nfs3\_max\_transfer\_size\_cots**

Description	Controls the maximum size of the data portion of an NFS version 3 READ, WRITE, REaddir, or REaddirPLUS request over TCP. This parameter controls both the maximum size of the request that the server returns as well as the maximum size of the request that the client generates.
Data Type	Integer (32-bit)
Default	1048576 bytes
Range	0 to $2^{31} - 1$
Units	Bytes
Dynamic?	Yes, but this parameter is set per file system at mount time. To affect a particular file system, unmount and mount the file system after changing this parameter.
Validation	None. However, setting the maximum transfer size on the server to 0 is likely to cause clients to malfunction or just decide not to attempt to talk to the server.
When to Change	Do not change this parameter unless transfer sizes larger than 1 Mbyte are preferred.
Commitment Level	Unstable

## **nfsrv Module Parameters**

This section describes NFS parameters for the `nfsrv` module.

### **nfsrv:nfs\_portmon**

Description	Controls some security checking that the NFS server attempts to do to enforce integrity on the part of its clients. The NFS server can check
-------------	--

whether the source port from which a request was sent was a *reserved port*. A reserved port has a number less than 1024. For BSD-based systems, these ports are reserved for processes being run by root. This security checking can prevent users from writing their own RPC-based applications that defeat the access checking that the NFS client uses.

Data Type	Integer (32-bit)
Default	0 (security checking disabled)
Range	0 (security checking disabled) or 1 (security checking enabled)
Units	Boolean values
Dynamic?	Yes
Validation	None
When to Change	Use this parameter to prevent malicious users from gaining access to files by using the NFS server that they would not ordinarily have access to. However, the <i>reserved port</i> notion is not universally supported. Thus, the security aspects of the check are very weak. Also, not all NFS client implementations bind their transport endpoints to a port number in the reserved range. Thus, interoperability problems might result if the security checking is enabled.
Commitment Level	Unstable

## nfssrv:rfs\_write\_async

**Description** Controls the behavior of the NFS version 2 server when it processes WRITE requests. The NFS version 2 protocol mandates that all modified data and metadata associated with the WRITE request reside on stable storage before the server can respond to the client. NFS version 2 WRITE requests are limited to 8192 bytes of data. Thus, each WRITE request might cause multiple small writes to the storage subsystem. This can cause a performance problem.

One method to accelerate NFS version 2 WRITE requests is to take advantage of a client behavior. Clients tend to send WRITE requests in batches. The server can take advantage of this behavior by clustering together the different WRITE requests into a single request to the underlying file system. Thus, the data to be written to the storage subsystem can be written in fewer, larger requests. This method can significantly increase the throughput for WRITE requests.

Data Type	Integer (32-bit)
-----------	------------------

Default	1 (clustering enabled)
Range	0 (clustering disabled) or 1 (clustering enabled)
Units	Boolean values
Dynamic?	Yes
Validation	None
When to Change	Some very small NFS clients, particularly PC clients, might not batch WRITE requests. Thus, the behavior required from the clients might not exist. In addition, the clustering in the NFS version 2 server might just add overhead and slow down performance instead of increasing it.
Commitment Level	Unstable

## **nfsrv:nfsauth\_ch\_cache\_max**

Description	Controls the size of the cache of client handles that contact the NFS authentication server. This server authenticates NFS clients to determine whether they are allowed access to the file handle that they are trying to use.
Data Type	Integer (32-bit)
Default	16
Range	0 to $2^{31} - 1$
Units	Client handles
Dynamic?	Yes
Validation	None
When to Change	This cache is not dynamic, so attempts to allocate a client handle when all are busy will fail. This failure results in requests being dropped by the NFS server because they could not be authenticated. Most often, this result is not a problem because the NFS client just times out and retransmits the request. However, for soft-mounted file systems on the client, the client might time out, not retry the request, and then return an error to the application. This situation might be avoided if you ensure that the size of the cache on the server is large enough to handle the load.
Commitment Level	Unstable

## **nfssrv:exi\_cache\_time**

Description	Controls the duration of time that entries are held in the NFS authentication cache before being purged due to memory pressure in the system.
Data Type	Long integer (32 bits on 32-bit platforms and 64 bits on 64-bit platforms)
Default	3600 seconds (1 hour)
Range	0 to $2^{31} - 1$ on 32-bit platforms 0 to $2^{63} - 1$ on 64-bit platforms
Units	Seconds
Dynamic?	Yes
Validation	None
When to Change	The size of the NFS authentication cache can be adjusted by varying the minimum age of entries that can get purged from the cache. The size of the cache should be controlled so that it is not allowed to grow too large, thus using system resources that are not allowed to be released due to this aging process.
Commitment Level	Unstable

## **rpcmod Module Parameters**

This section describes NFS parameters for the rpcmod module.

### **rpcmod:clnt\_max\_conns**

Description	Controls the number of TCP connections that the NFS client uses when communicating with each NFS server. The kernel RPC is constructed so that it can multiplex RPCs over a single connection. However, multiple connections can be used, if preferred.
Data Type	Integer (32-bit)
Default	1
Range	1 to $2^{31} - 1$
Units	Connections

Dynamic?	Yes
Validation	None
When to Change	In general, one connection is sufficient to achieve full network bandwidth. However, if TCP cannot utilize the bandwidth offered by the network in a single stream, then multiple connections might increase the throughput between the client and the server.  Increasing the number of connections doesn't come without consequences. Increasing the number of connections also increases kernel resource usage needed to keep track of each connection.
Commitment Level	Unstable

## rpcmod:clnt\_idle\_timeout

Description	Controls the duration of time on the client that a connection between the client and server is allowed to remain idle before being closed.
Data Type	Long integer (32 bits on 32-bit platforms and 64 bits on 64-bit platforms)
Default	300,000 milliseconds (5 minutes)
Range	0 to $2^{31} - 1$ on 32-bit platforms 0 to $2^{63} - 1$ on 64-bit platforms
Units	Milliseconds
Dynamic?	Yes
Validation	None
When to Change	Use this parameter to change the time that idle connections are allowed to exist on the client before being closed. You might want to close connections at a faster rate to avoid consuming system resources.
Commitment Level	Unstable

## rpcmod:svc\_idle\_timeout

Description	Controls the duration of time on the server that a connection between the client and server is allowed to remain idle before being closed.
Data Type	Long integer (32 bits on 32-bit platforms and 64 bits on 64-bit platforms)

Default	360,000 milliseconds (6 minutes)
Range	0 to $2^{31} - 1$ on 32-bit platforms 0 to $2^{63} - 1$ on 64-bit platforms
Units	Milliseconds
Dynamic?	Yes
Validation	None
When to Change	Use this parameter to change the time that idle connections are allowed to exist on the server before being closed. You might want to close connections at a faster rate to avoid consuming system resources.
Commitment Level	Unstable

## **rpcmod:svc\_default\_stksize**

Description	Sets the size of the kernel stack for kernel RPC service threads.
Data Type	Integer (32-bit)
Default	The default value is 0. This value means that the stack size is set to the system default.
Range	0 to $2^{31} - 1$
Units	Bytes
Dynamic?	Yes, for all new threads that are allocated. The stack size is set when the thread is created. Therefore, changes to this parameter do not affect existing threads but are applied to all new threads that are allocated.
Validation	None
When to Change	Very deep call depths can cause the stack to overflow and cause red zone faults. The combination of a fairly deep call depth for the transport, coupled with a deep call depth for the local file system, can cause NFS service threads to overflow their stacks.  Set this parameter to a multiple of the hardware page size on the platform.
Commitment Level	Unstable



## rpcmod:svc\_default\_max\_same\_xprt

Description	Controls the maximum number of requests that are processed for each transport endpoint before switching transport endpoints. The kernel RPC works by having a pool of service threads and a pool of transport endpoints. Any one of the service threads can process requests from any one of the transport endpoints. For performance, multiple requests on each transport endpoint are consumed before switching to a different transport endpoint. This approach offers performance benefits while avoiding starvation.
Data Type	Integer (32-bit)
Default	8
Range	0 to $2^{31} - 1$
Units	Requests
Dynamic?	Yes, but the maximum number of requests to process before switching transport endpoints is set when the transport endpoint is configured into the kernel RPC subsystem. Changes to this parameter only affect new transport endpoints, not existing transport endpoints.
Validation	None
When to Change	Tune this parameter so that services can take advantage of client behaviors such as the clustering that accelerate NFS version 2 WRITE requests. Increasing this parameter might result in the server being better able to take advantage of client behaviors.
Commitment Level	Unstable

## rpcmod:maxdupreqs

Description	Controls the size of the duplicate request cache that detects RPC- level retransmissions on connectionless transports. This cache is indexed by the client network address and the RPC procedure number, program number, version number, and transaction ID. This cache avoids processing retransmitted requests that might not be idempotent.
Data Type	Integer (32-bit)
Default	1024
Range	1 to $2^{31} - 1$
Units	Requests

Dynamic?	<p>The cache is dynamically sized, but the hash queues that provide fast access to the cache are statically sized. Making the cache very large might result in long search times to find entries in the cache.</p> <p>Do not set the value of this parameter to 0. This value prevents the NFS server from handling non idempotent requests.</p>
Validation	None
When to Change	<p>Examine the value of this parameter if false failures are encountered by NFS clients. For example, if an attempt to create a directory fails, but the directory is actually created, perhaps that retransmitted MKDIR request was not detected by the server.</p> <p>The size of the cache should match the load on the server. The cache records non idempotent requests and so only needs to track a portion of the total requests. The cache does need to hold the information long enough to be able to detect a retransmission by the client. Typically, the client timeout for connectionless transports is relatively short, starting around 1 second and increasing to about 20 seconds.</p>
Commitment Level	Unstable

## **rpcmod:cotsmaxdupreqs**

Description	<p>Controls the size of the duplicate request cache that detects RPC- level retransmissions on connection-oriented transports. This cache is indexed by the client network address and the RPC procedure number, program number, version number, and transaction ID. This cache avoids processing retransmitted requests that might not be idempotent.</p>
Data Type	Integer (32-bit)
Default	1024
Range	1 to $2^{31} - 1$
Units	Requests
Dynamic?	Yes
Validation	<p>The cache is dynamically sized, but the hash queues that provide fast access to the cache are statically sized. Making the cache very large might result in long search times to find entries in the cache.</p> <p>Do not set the value of this parameter to 0. It prevents the NFS server from handling non-idempotent requests.</p>

When to Change	<p>Examine the value of this parameter if false failures are encountered by NFS clients. For example, if an attempt to create a directory fails, but the directory is actually created, it is possible that a retransmitted MKDIR request was not detected by the server.</p> <p>The size of the cache should match the load on the server. The cache records non-idempotent requests and so only needs to track a portion of the total requests. It does need to hold the information long enough to be able to detect a retransmission on the part of the client. Typically, the client timeout for connection oriented transports is very long, about 1 minute. Thus, entries need to stay in the cache for fairly long times.</p>
Commitment Level	Unstable



# Internet Protocol Suite Tunable Parameters

---

This chapter describes various Internet Protocol suite parameters or properties.

- “IP Tunable Parameters” on page 134
- “TCP Tunable Parameters” on page 136
- “UDP Tunable Parameters” on page 151
- “IPQoS Tunable Parameter” on page 152
- “SCTP Tunable Parameters” on page 153
- “Per-Route Metrics” on page 160

## Where to Find Tunable Parameter Information

Tunable Parameter	For Information
Solaris kernel tunables	Chapter 2, “Oracle Solaris Kernel Tunable Parameters”
NFS tunable parameters	Chapter 3, “NFS Tunable Parameters”
Network Cache and Accelerator (NCA) tunable parameters	Chapter 5, “Network Cache and Accelerator Tunable Parameters”

## Overview of Tuning IP Suite Parameters

You can set all of the tuning parameters described in this chapter by using the `ipadm` command except for the following parameters:

- “`ipcl_conn_hash_size`” on page 146
- “`ip_squeue_worker_wait`” on page 147
- “`ip_squeue_fanout`” on page 136

These parameters can only be set in the `/etc/system` file.

Use the following syntax to set TCP/IP parameters by using the `ipadm` command:

```
# ipadm set-prop -p parameter tcp|ip
```

For example:

```
# ipadm set-prop -p extra_priv_ports=1047 tcp
# ipadm show-prop -p extra_priv_ports tcp
PROTO PROPERTY          PERM CURRENT    PERSISTENT  DEFAULT    POSSIBLE
tcp  extra_priv_ports    rw   1047          1047        2049,4045  1-65535
```

For more information, see [ipadm\(1M\)](#).

Use the following syntax to set TCP/IP parameters by using the `ndd` command:

```
# ndd -set driver parameter
```

For more information, see [ndd\(1M\)](#).

Although the SMF framework provides a method for managing system services, `ipadm` commands are still included in system startup scripts. For more information on creating a startup script, see “Using Run Control Scripts” in *System Administration Guide: Basic Administration*.

## IP Suite Parameter Validation

All parameters described in this section are checked to verify that they fall in the parameter range. The parameter's range is provided with the description for each parameter.

## Internet Request for Comments (RFCs)

Internet protocol and standard specifications are described in RFC documents. You can get copies of RFCs from `ftp://ftp.rfc-editor.org/in-notes`. Browse RFC topics by viewing the `rfc-index.txt` file at this site.

# IP Tunable Parameters

## `_icmp_err_interval` and `_icmp_err_burst`

Description	Controls the rate of IP in generating IPv4 ICMP error messages. IP generates only up to <code>_icmp_err_burst</code> IPv4 error messages in any <code>_icmp_err_interval</code> .
-------------	---

The `_icmp_err_interval` parameter protects IP from denial of service attacks. Setting this parameter to 0 disables rate limiting. It does not disable the generation of error messages.

Default	100 milliseconds for <code>_icmp_err_interval</code> 10 error messages for <code>ip_icmp_err_burst</code>
Range	0 – 99,999 milliseconds for <code>_icmp_err_interval</code> 1 – 99,999 error messages for <code>_icmp_err_burst</code>
Dynamic?	Yes
When to Change	If you need a higher error message generation rate for diagnostic purposes.
Commitment Level	Unstable

## **`_respond_to_echo_broadcast`**

Description	Controls whether IPv4 responds to a broadcast ICMPv4 echo request.
Default	1 (enabled)
Range	0 (disabled) or 1 (enabled)
Dynamic?	Yes
When to Change	If you do not want this behavior for security reasons, disable it.
Commitment Level	Unstable

## **`_addrs_per_if`**

Description	Defines the maximum number of logical interfaces associated with a real interface.
Default	256
Range	1 to 8192
Dynamic?	Yes
When to Change	Do not change the value. If more logical interfaces are required, you might consider increasing the value. However, recognize that this change might have a negative impact on IP's performance.
Commitment Level	Unstable

## ip\_queue\_fanout

Description	Determines the mode of associating TCP/IP connections with queues  A value of 0 associates a new TCP/IP connection with the CPU that creates the connection. A value of 1 associates the connection with multiple queues that belong to different CPUs.
Default	0
Range	0 or 1
Dynamic?	Yes
When to Change	Consider setting this parameter to 1 to spread the load across all CPUs in certain situations. For example, when the number of CPUs exceed the number of NICs, and one CPU is not capable of handling the network load of a single NIC, change this parameter to 1.  This property can only be set in the <code>/etc/system</code> file.
Zone Configuration	This parameter can only be set in the global zone.
Commitment Level	Unstable

## TCP Tunable Parameters

### \_deferred\_ack\_interval

Description	Specifies the time-out value for the TCP-delayed acknowledgment (ACK) timer for hosts that are not directly connected.  Refer to RFC 1122, 4.2.3.2.
Default	100 milliseconds
Range	1 millisecond to 1 minute
Dynamic?	Yes
When to Change	Do not increase this value to more than 500 milliseconds.  Increase the value under the following circumstances: <ul style="list-style-type: none"><li>▪ Slow network links (less than 57.6 Kbps) with greater than 512 bytes maximum segment size (MSS)</li><li>▪ The interval for receiving more than one TCP segment is short</li></ul>



Commitment Level Unstable

## **`_local_dack_interval`**

Description	Specifies the time-out value for TCP-delayed acknowledgment (ACK) timer for hosts that are directly connected.  Refer to RFC 1122, 4.2.3.2.
Default	50 milliseconds
Range	10 milliseconds to 500 milliseconds
Dynamic?	Yes
When to Change	Do not increase this value to more than 500 milliseconds.  Increase the value under the following circumstances: <ul style="list-style-type: none"> <li>▪ Slow network links (less than 57.6 Kbps) with greater than 512 bytes maximum segment size (MSS)</li> <li>▪ The interval for receiving more than one TCP segment is short</li> </ul>
Commitment Level	Unstable

## **`_deferred_acks_max`**

Description	Specifies the maximum number of TCP segments received from remote destinations (not directly connected) before an acknowledgment (ACK) is generated. TCP segments are measured in units of maximum segment size (MSS) for individual connections. If set to 0 or 1, no ACKs are delayed, assuming all segments are 1 MSS long. The actual number is dynamically calculated for each connection. The value is the default maximum.
Default	2
Range	0 to 16
Dynamic?	Yes
When to Change	Do not change the value. In some circumstances, when the network traffic becomes very bursty because of the delayed ACK effect, decrease the value. Do not decrease this value below 2.
Commitment Level	Unstable

## **`_local_dacks_max`**

Description	Specifies the maximum number of TCP segments received from directly connected destinations before an acknowledgment (ACK) is generated. TCP segments are measured in units of maximum segment size (MSS) for individual connections. If set to 0 or 1, it means no ACKs are delayed, assuming all segments are 1 MSS long. The actual number is dynamically calculated for each connection. The value is the default maximum.
Default	8
Range	0 to 16
Dynamic?	Yes
When to Change	Do not change the value. In some circumstances, when the network traffic becomes very bursty because of the delayed ACK effect, decrease the value. Do not decrease this value below 2.
Commitment Level	Unstable

## **`_wscale_always`**

Description	When this parameter is enabled, which is the default setting, TCP always sends a SYN segment with the window scale option, even if the window scale option value is 0. Note that if TCP receives a SYN segment with the window scale option, even if the parameter is disabled, TCP responds with a SYN segment with the window scale option. In addition, the option value is set according to the receive window size.  Refer to RFC 1323 for the window scale option.
Default	1 (enabled)
Range	0 (disabled) or 1 (enabled)
Dynamic?	Yes
When to Change	If there is an interoperability problem with an old TCP stack that does not support the window scale option, disable this parameter.
Commitment Level	Unstable

## **\_tstamp\_always**

Description	If set to 1, TCP always sends a SYN segment with the timestamp option. Note that if TCP receives a SYN segment with the timestamp option, TCP responds with a SYN segment with the timestamp option even if the parameter is set to 0.
Default	0 (disabled)
Range	0 (disabled) or 1 (enabled)
Dynamic?	Yes
When to Change	If getting an accurate measurement of round-trip time (RTT) and TCP sequence number wraparound is a problem, enable this parameter.  Refer to RFC 1323 for more reasons to enable this option.
Commitment Level	Unstable

## **send\_maxbuf**

Description	Defines the default send window size in bytes. Refer to “ <a href="#">Per-Route Metrics</a> ” on page 160 for a discussion of setting a different value on a per-route basis. See also “ <a href="#">_max_buf</a> ” on page 140.
Default	49,152
Range	4096 to 1,073,741,824
Dynamic?	Yes
When to Change	An application can use <code>setsockopt(3XNET) SO_SNDBUF</code> to change the individual connection's send buffer.
Commitment Level	Unstable

## **recv\_maxbuf**

Description	Defines the default receive window size in bytes. Refer to “ <a href="#">Per-Route Metrics</a> ” on page 160 for a discussion of setting a different value on a per-route basis. See also “ <a href="#">_max_buf</a> ” on page 140 and “ <a href="#">_recv_hiwat_minmss</a> ” on page 151.
Default	128,000
Range	2048 to 1,073,741,824

Dynamic?	Yes
When to Change	An application can use <code>setsockopt(3XNET)</code> <code>SO_RCVBUF</code> to change the individual connection's receive buffer.
Commitment Level	Unstable

## **`_max_buf`**

Description	Defines the maximum buffer size in bytes. This parameter controls how large the send and receive buffers are set to by an application that uses <code>setsockopt(3XNET)</code> .
Default	1,048,576
Range	8192 to 1,073,741,824
Dynamic?	Yes
When to Change	If TCP connections are being made in a high-speed network environment, increase the value to match the network link speed.
Commitment Level	Unstable

## **`_cwnd_max`**

Description	Defines the maximum value of the TCP congestion window (cwnd) in bytes.  For more information on the TCP congestion window, refer to RFC 1122 and RFC 2581.
Default	1,048,576
Range	128 to 1,073,741,824
Dynamic?	Yes
When to Change	Even if an application uses <code>setsockopt(3XNET)</code> to change the window size to a value higher than <code>_cwnd_max</code> , the actual window used can never grow beyond <code>_cwnd_max</code> . Thus, <code>_max_buf</code> should be greater than <code>_cwnd_max</code> .
Commitment Level	Unstable

## **`_slow_start_initial`**

Description	<p>Defines the maximum initial congestion window (cwnd) size in the maximum segment size (MSS) of a TCP connection.</p> <p>Refer to RFC 2414 on how the initial congestion window size is calculated.</p>
Default	4
Range	1 to 4
Dynamic?	Yes
When to Change	<p>Do not change the value.</p> <p>If the initial cwnd size causes network congestion under special circumstances, decrease the value.</p>
Commitment Level	Unstable

## **`_slow_start_after_idle`**

Description	<p>The congestion window size in the maximum segment size (MSS) of a TCP connection after it has been idled (no segment received) for a period of one retransmission timeout (RTO).</p> <p>Refer to RFC 2414 on how the initial congestion window size is calculated.</p>
Default	4
Range	1 to 16,384
Dynamic?	Yes
When to Change	For more information, see “ <a href="#">_slow_start_initial</a> ” on page 141.
Commitment Level	Unstable

## **sack**

Description	<p>If set to 2, TCP always sends a SYN segment with the selective acknowledgment (SACK) permitted option. If TCP receives a SYN segment with a SACK-permitted option and this parameter is set to 1, TCP responds with a SACK-permitted option. If the parameter is set to</p>
-------------	--

0, TCP does not send a SACK-permitted option, regardless of whether the incoming segment contains the SACK permitted option.

Refer to RFC 2018 for information on the SACK option.

Default	2 (active enabled)
Range	0 (disabled), 1 (passive enabled), or 2 (active enabled)
Dynamic?	Yes
When to Change	SACK processing can improve TCP retransmission performance so it should be actively enabled. Sometimes, the other side can be confused with the SACK option actively enabled. If this confusion occurs, set the value to 1 so that SACK processing is enabled only when incoming connections allow SACK processing.
Commitment Level	Unstable

## **`_rev_src_routes`**

Description	If set to 0, TCP does not reverse the IP source routing option for incoming connections for security reasons. If set to 1, TCP does the normal reverse source routing.
Default	0 (disabled)
Range	0 (disabled) or 1 (enabled)
Dynamic?	Yes
When to Change	If IP source routing is needed for diagnostic purposes, enable it.
Commitment Level	Unstable

## **`_time_wait_interval`**

Description	Specifies the time in milliseconds that a TCP connection stays in TIME-WAIT state.  For more information, refer to RFC 1122, 4.2.2.13.
Default	60,000 (60 seconds)
Range	1 second to 10 minutes
Dynamic?	Yes
When to Change	Do not set the value lower than 60 seconds.

	For information on changing this parameter, refer to RFC 1122, 4.2.2.13.
Commitment Level	Unstable

## ecn

Description	<p>Controls Explicit Congestion Notification (ECN) support.</p> <p>If this parameter is set to 0, TCP does not negotiate with a peer that supports the ECN mechanism.</p> <p>If this parameter is set to 1 when initiating a connection, TCP does not tell a peer that it supports ECN mechanism.</p> <p>However, TCP tells a peer that it supports ECN mechanism when accepting a new incoming connection request if the peer indicates that it supports ECN mechanism in the SYN segment.</p> <p>If this parameter is set to 2, in addition to negotiating with a peer on the ECN mechanism when accepting connections, TCP indicates in the outgoing SYN segment that it supports the ECN mechanism when TCP makes active outgoing connections.</p> <p>Refer to RFC 3168 for information on ECN.</p>
Default	1 (passive enabled)
Range	0 (disabled), 1 (passive enabled), or 2 (active enabled)
Dynamic?	Yes
When to Change	<p>ECN can help TCP better handle congestion control. However, there are existing TCP implementations, firewalls, NATs, and other network devices that are confused by this mechanism. These devices do not comply to the IETF standard.</p> <p>Because of these devices, the default value of this parameter is set to 1. In rare cases, passive enabling can still cause problems. Set the parameter to 0 only if absolutely necessary.</p>
Commitment Level	Unstable

## **`_conn_req_max_q`**

Description	Specifies the default maximum number of pending TCP connections for a TCP listener waiting to be accepted by <code>accept(3SOCKET)</code> . See also “ <code>_conn_req_max_q0</code> ” on page 144.
Default	128
Range	1 to 4,294,967,295
Dynamic?	Yes
When to Change	<p>For applications such as web servers that might receive several connection requests, the default value might be increased to match the incoming rate.</p> <p>Do not increase the parameter to a very large value. The pending TCP connections can consume excessive memory. Also, if an application cannot handle that many connection requests fast enough because the number of pending TCP connections is too large, new incoming requests might be denied.</p> <p>Note that increasing <code>_conn_req_max_q</code> does not mean that applications can have that many pending TCP connections. Applications can use <code>listen(3SOCKET)</code> to change the maximum number of pending TCP connections for each socket. This parameter is the maximum an application can use <code>listen()</code> to set the number to. Thus, even if this parameter is set to a very large value, the actual maximum number for a socket might be much less than <code>_conn_req_max_q</code>, depending on the value used in <code>listen()</code>.</p>
Commitment Level	Unstable

## **`_conn_req_max_q0`**

Description	<p>Specifies the default maximum number of incomplete (three-way handshake not yet finished) pending TCP connections for a TCP listener.</p> <p>For more information on TCP three-way handshake, refer to RFC 793. See also “<code>_conn_req_max_q</code>” on page 144.</p>
Default	1024
Range	0 to 4,294,967,296
Dynamic?	Yes



**When to Change** For applications such as web servers that might receive excessive connection requests, you can increase the default value to match the incoming rate.

The following explains the relationship between `_conn_req_max_q0` and the maximum number of pending connections for each socket.

When a connection request is received, TCP first checks if the number of pending TCP connections (three-way handshake is done) waiting to be accepted exceeds the maximum ( $N$ ) for the listener. If the connections are excessive, the request is denied. If the number of connections is allowable, then TCP checks if the number of incomplete pending TCP connections exceeds the sum of  $N$  and `_conn_req_max_q0`. If it does not, the request is accepted. Otherwise, the oldest incomplete pending TCP request is dropped.

**Commitment Level** Unstable

## `_conn_req_min`

**Description** Specifies the default minimum value for the maximum number of pending TCP connection requests for a listener waiting to be accepted. This is the lowest maximum value of `listen(3SOCKET)` that an application can use.

**Default** 1

**Range** 1 to 1024

**Dynamic?** Yes

**When to Change** This parameter can be a solution for applications that use `listen(3SOCKET)` to set the maximum number of pending TCP connections to a value too low. Increase the value to match the incoming connection request rate.

**Commitment Level** Unstable

## `_rst_sent_rate_enabled`

**Description** If this parameter is set to 1, the maximum rate of sending a RST segment is controlled by the `ipadm` parameter, `_rst_sent_rate`. If this parameter is set to 0, no rate control when sending a RST segment is available.

Default	1 (enabled)
Range	0 (disabled) or 1 (enabled)
Dynamic?	Yes
When to Change	This tunable helps defend against denial of service attacks on TCP by limiting the rate by which a RST segment is sent out. The only time this rate control should be disabled is when strict conformance to RFC 793 is required.
Commitment Level	Unstable

## **`_rst_sent_rate`**

Description	Sets the maximum number of RST segments that TCP can send out per second.
Default	40
Range	0 to 4,294,967,295
Dynamic?	Yes
When to Change	In a TCP environment, there might be a legitimate reason to generate more RSTs than the default value allows. In this case, increase the default value of this parameter.
Commitment Level	Unstable

## **TCP/IP Parameters Set in the `/etc/system` File**

The following parameters can be set only in the `/etc/system` file. After the file is modified, reboot the system.

For example, the following entry sets the `ipcl_conn_hash_size` parameter:

```
set ip:ipcl_conn_hash_sizes=value
```

### **`ipcl_conn_hash_size`**

Description	Controls the size of the connection hash table used by IP. The default value of 0 means that the system automatically sizes an appropriate value for this parameter at boot time, depending on the available memory.
Data Type	Unsigned integer

Default	0
Range	0 to 82,500
Dynamic?	No. The parameter can only be changed at boot time.
When to Change	If the system consistently has tens of thousands of TCP connections, the value can be increased accordingly. Increasing the hash table size means that more memory is wired down, thereby reducing available memory to user applications.
Commitment Level	Unstable

### **ip\_squeue\_worker\_wait**

Description	Governs the maximum delay in waking up a worker thread to process TCP/IP packets that are enqueued on an squeue. An <i>squeue</i> is a serialization queue that is used by the TCP/IP kernel code to process TCP/IP packets.
Default	10 milliseconds
Range	0 – 50 milliseconds
Dynamic?	Yes
When to Change	Consider tuning this parameter if latency is an issue, and network traffic is light. For example, if the machine serves mostly interactive network traffic.  The default value usually works best on a network file server, a web server, or any server that has substantial network traffic.
Zone Configuration	This parameter can only be set in the global zone.
Commitment Level	Unstable

## **TCP Parameters With Additional Cautions**

Changing the following parameters is not recommended.

### **\_keepalive\_interval**

Description	This <code>ipadm</code> parameter sets a probe interval that is first sent out after a TCP connection is idle on a system-wide basis.
-------------	---

Solaris supports the TCP keep-alive mechanism as described in RFC 1122. This mechanism is enabled by setting the `SO_KEEPALIVE` socket option on a TCP socket.

If `SO_KEEPALIVE` is enabled for a socket, the first keep-alive probe is sent out after a TCP connection is idle for two hours, the default value of the `tcp_keepalive_interval` parameter. If the peer does not respond to the probe after eight minutes, the TCP connection is aborted. For more information, refer to “[\\_rexmit\\_interval\\_initial](#)” on page 149.

You can also use the `TCP_KEEPALIVE_THRESHOLD` socket option on individual applications to override the default interval so that each application can have its own interval on each socket. The option value is an unsigned integer in milliseconds. See also [tcp\(7P\)](#).

Default	2 hours
Range	10 seconds to 10 days
Units	Unsigned integer (milliseconds)
Dynamic?	Yes
When to Change	Do not change the value. Lowering it may cause unnecessary network traffic and might also increase the chance of premature termination of the connection because of a transient network problem.
Commitment Level	Unstable

### **\_ip\_abort\_interval**

**Description** Specifies the default total retransmission timeout value for a TCP connection. For a given TCP connection, if TCP has been retransmitting for `_ip_abort_interval` period of time and it has not received any acknowledgment from the other endpoint during this period, TCP closes this connection.

For TCP retransmission timeout (RTO) calculation, refer to RFC 1122, 4.2.3. See also “[\\_rexmit\\_interval\\_max](#)” on page 149.

Default	8 minutes
Range	500 milliseconds to 1193 hours
Dynamic?	Yes
When to Change	Do not change this value. See “ <a href="#">_rexmit_interval_max</a> ” on page 149 for exceptions.

Commitment Level Unstable

### **`_rexmit_interval_initial`**

Description	Specifies the default initial retransmission timeout (RTO) value for a TCP connection. Refer to <a href="#">“Per-Route Metrics” on page 160</a> for a discussion of setting a different value on a per-route basis.
Default	3 seconds
Range	1 millisecond to 20 seconds
Dynamic?	Yes
When to Change	Do not change this value. Lowering the value can result in unnecessary retransmissions.
Commitment Level	Unstable

### **`_rexmit_interval_max`**

Description	Defines the default maximum retransmission timeout value (RTO). The calculated RTO for all TCP connections cannot exceed this value. See also <a href="#">“_ip_abort_interval” on page 148</a> .
Default	60 seconds
Range	1 millisecond to 2 hours
Dynamic?	Yes
When to Change	Do not change the value in a normal network environment.  If, in some special circumstances, the round-trip time (RTT) for a connection is about 10 seconds, you can increase this value. If you change this value, you should also change the <code>_ip_abort_interval</code> parameter. Change the value of <code>_ip_abort_interval</code> to at least four times the value of <code>_rexmit_interval_max</code> .
Commitment Level	Unstable

### **`_rexmit_interval_min`**

Description	Specifies the default minimum retransmission time out (RTO) value. The calculated RTO for all TCP connections cannot be lower than this value. See also <a href="#">“_rexmit_interval_max” on page 149</a> .
Default	400 milliseconds
Range	1 millisecond to 20 seconds

Dynamic?	Yes
When to Change	Do not change the value in a normal network environment.  TCP's RTO calculation should cope with most RTT fluctuations. If, in some very special circumstances, the round-trip time (RTT) for a connection is about 10 seconds, increase this value. If you change this value, you should change the <code>_rexmit_interval_max</code> parameter. Change the value of <code>_rexmit_interval_max</code> to at least eight times the value of <code>_rexmit_interval_min</code> .
Commitment Level	Unstable

### **`_rexmit_interval_extra`**

Description	Specifies a constant added to the calculated retransmission time out value (RTO).
Default	0 milliseconds
Range	0 to 2 hours
Dynamic?	Yes
When to Change	Do not change the value.  When the RTO calculation fails to obtain a good value for a connection, you can change this value to avoid unnecessary retransmissions.
Commitment Level	Unstable

### **`_tstamp_if_wscale`**

Description	If this parameter is set to 1, and the window scale option is enabled for a connection, TCP also enables the <code>timestamp</code> option for that connection.
Default	1 (enabled)
Range	0 (disabled) or 1 (enabled)
Dynamic?	Yes
When to Change	Do not change this value. In general, when TCP is used in high-speed network, protection against sequence number wraparound is essential. Thus, you need the <code>timestamp</code> option.
Commitment Level	Unstable

**`_recv_hiwat_minmss`**

Description	Controls the default minimum receive window size. The minimum is <code>_recv_hiwat_minmss</code> times the size of maximum segment size (MSS) of a connection.
Default	8
Range	1 to 65,536
Dynamic?	Yes
When to Change	Do not change the value. If changing it is necessary, do not change the value lower than 4.
Commitment Level	Unstable

## UDP Tunable Parameters

**`send_maxbuf`**

Description	Defines the default maximum UDP socket datagram size. .
Default	57,344 bytes
Range	1,024 to 1,073,741,824 bytes
Dynamic?	Yes
When to Change	Note that an application can use <code>setsockopt(3XNET)</code> <code>SO_SNDBUF</code> to change the size for an individual socket. In general, you do not need to change the default value.
Commitment Level	Unstable

**`recv_maxbuf`**

Description	Defines the default maximum UDP socket receive buffer size.
Default	57,344 bytes
Range	128 to 1,073,741,824 bytes
Dynamic?	Yes

When to Change Note that an application can use `setsockopt(3XNET)` `SO_RCVBUF` to change the size for an individual socket. In general, you do not need to change the default value.

Commitment Level Unstable

## IPQoS Tunable Parameter

### `_policy_mask`

Description Enables or disables IPQoS processing in any of the following callout positions: forward outbound, forward inbound, local outbound, and local inbound. This parameter is a bitmask as follows:

Not Used	Not Used	Not Used	Not Used	Forward Outbound	Forward Inbound	Local Outbound	Local Inbound
X	X	X	X	0	0	0	0

A 1 in any of the position masks or disables IPQoS processing in that particular callout position. For example, a value of `0x01` disables IPQoS processing for all the local inbound packets.

Default The default value is 0, meaning that IPQoS processing is enabled in all the callout positions.

Range 0 (0x00) to 15 (0x0F). A value of 15 indicates that IPQoS processing is disabled in all the callout positions.

Dynamic? Yes

When to Change If you want to enable or disable IPQoS processing in any of the callout positions.

Commitment Level Unstable



## SCTP Tunable Parameters

### **`_max_init_retr`**

Description	Controls the maximum number of attempts an SCTP endpoint should make at resending an INIT chunk. The SCTP endpoint can use the SCTP initiation structure to override this value.
Default	8
Range	0 to 128
Dynamic?	Yes
When to Change	The number of INIT retransmissions depend on “ <a href="#">_pa_max_retr</a> ” on <a href="#">page 153</a> . Ideally, <code>_max_init_retr</code> should be less than or equal to <code>_pa_max_retr</code> .
Commitment Level	Unstable

### **`_pa_max_retr`**

Description	Controls the maximum number of retransmissions (over all paths) for an SCTP association. The SCTP association is aborted when this number is exceeded.
Default	10
Range	1 to 128
Dynamic?	Yes
When to Change	The maximum number of retransmissions over all paths depend on the number of paths and the maximum number of retransmission over each path. Ideally, <code>sctp_pa_max_retr</code> should be set to the sum of “ <a href="#">_pp_max_retr</a> ” on <a href="#">page 154</a> over all available paths. For example, if there are 3 paths to the destination and the maximum number of retransmissions over each of the 3 paths is 5, then <code>_pa_max_retr</code> should be set to less than or equal to 15. (See the Note in Section 8.2, RFC 2960.)
Commitment Level	Unstable

## **\_pp\_max\_retr**

Description	Controls the maximum number of retransmissions over a specific path. When this number is exceeded for a path, the path (destination) is considered unreachable.
Default	5
Range	1 to 128
Dynamic?	Yes
When to Change	Do not change this value to less than 5.
Commitment Level	Unstable

## **\_cwnd\_max**

Description	Controls the maximum value of the congestion window for an SCTP association.
Default	1,048,576
Range	128 to 1,073,741,824
Dynamic?	Yes
When to Change	Even if an application uses <code>setsockopt(3XNET)</code> to change the window size to a value higher than <code>_cwnd_max</code> , the actual window used can never grow beyond <code>_cwnd_max</code> . Thus, “ <code>_max_buf</code> ” on page 157 should be greater than <code>_cwnd_max</code> .
Commitment Level	Unstable

## **\_ipv4\_ttl**

Description	Controls the time to live (TTL) value in the IP version 4 header for the outbound IP version 4 packets on an SCTP association.
Default	64
Range	1 to 255
Dynamic?	Yes
When to Change	Generally, you do not need to change this value. Consider increasing this parameter if the path to the destination is likely to span more than 64 hops.
Commitment Level	Unstable

## **`_heartbeat_interval`**

Description	Computes the interval between HEARTBEAT chunks to an idle destination, that is allowed to heartbeat.  An SCTP endpoint periodically sends an HEARTBEAT chunk to monitor the reachability of the idle destinations transport addresses of its peer.
Default	30 seconds
Range	0 to 86,400 seconds
Dynamic?	Yes
When to Change	Refer to RFC 2960, section 8.3.
Commitment Level	Unstable

## **`_new_secret_interval`**

Description	Determines when a new secret needs to be generated. The generated secret is used to compute the MAC for a cookie.
Default	2 minutes
Range	0 to 1,440 minutes
Dynamic?	Yes
When to Change	Refer to RFC 2960, section 5.1.3.
Commitment Level	Unstable

## **`_initial_mtu`**

Description	Determines the initial maximum send size for an SCTP packet including the length of the IP header.
Default	1500 bytes
Range	68 to 65,535
Dynamic?	Yes
When to Change	Increase this parameter if the underlying link supports frame sizes that are greater than 1500 bytes.
Commitment Level	Unstable

## **`_deferred_ack_interval`**

Description	Sets the time-out value for SCTP delayed acknowledgment (ACK) timer in milliseconds.
Default	100 milliseconds
Range	1 to 60,000 milliseconds
Dynamic?	Yes
When to Change	Refer to RFC 2960, section 6.2.
Commitment Level	Unstable

## **`_ignore_path_mtu`**

Description	Enables or disables path MTU discovery.
Default	0 (disabled)
Range	0 (disabled) or 1 (enabled)
Dynamic?	Yes
When to Change	Enable this parameter if you want to ignore MTU changes along the path. However, doing so might result in IP fragmentation if the path MTU decreases.
Commitment Level	Unstable

## **`_initial_ssthresh`**

Description	Sets the initial slow start threshold for a destination address of the peer.
Default	102,400
Range	1024 to 4,294,967,295
Dynamic?	Yes
When to Change	Refer to RFC 2960, section 7.2.1.
Commitment Level	Unstable

## **`_max_buf`**

Description	Controls the maximum buffer size in bytes. It controls how large the send and receive buffers are set to by an application that uses <code>getsockopt(3SOCKET)</code> .
Default	1,048,576
Range	8,192 to 1,073,741,824
Dynamic?	Yes
When to Change	Increase the value of this parameter to match the network link speed if associations are being made in a high-speed network environment.
Commitment Level	Unstable

## **`_ipv6_hoplimit`**

Description	Sets the value of the hop limit in the IP version 6 header for the outbound IP version 6 packets on an SCTP association.
Default	60
Range	0 to 255
Dynamic?	Yes
When to Change	Generally, you do not need to change this value. Consider increasing this parameter if the path to the destination is likely to span more than 60 hops.
Commitment Level	Unstable

## **`_rto_min`**

Description	Sets the lower bound for the retransmission timeout (RTO) in milliseconds for all the destination addresses of the peer.
Default	1,000
Range	500 to 60,000
Dynamic?	Yes
When to Change	Refer to RFC 2960, section 6.3.1.
Commitment Level	Unstable

## **`_rto_max`**

Description	Controls the upper bound for the retransmission timeout (RTO) in milliseconds for all the destination addresses of the peer.
Default	60,000
Range	1,000 to 60,000,000
Dynamic?	Yes
When to Change	Refer to RFC 2960, section 6.3.1.
Commitment Level	Unstable

## **`_rto_initial`**

Description	Controls the initial retransmission timeout (RTO) in milliseconds for all the destination addresses of the peer.
Default	3,000
Range	1,000 to 60,000,000
Dynamic?	Yes
When to Change	Refer to RFC 2960, section 6.3.1.
Commitment Level	Unstable

## **`_cookie_life`**

Description	Sets the lifespan of a cookie in milliseconds.
Default	60,000
Range	10 to 60,000,000
Dynamic?	Yes
When to Change	Generally, you do not need to change this value. This parameter might be changed in accordance with “ <a href="#">_rto_max</a> ” on page 158.
Commitment Level	Unstable

## **`_max_in_streams`**

Description	Controls the maximum number of inbound streams permitted for an SCTP association.
-------------	---

Default	32
Range	1 to 65,535
Dynamic?	Yes
When to Change	Refer to RFC 2960, section 5.1.1.
Commitment Level	Unstable

## **`_initial_out_streams`**

Description	Controls the maximum number of outbound streams permitted for an SCTP association.
Default	32
Range	1 to 65,535
Dynamic?	Yes
When to Change	Refer to RFC 2960, section 5.1.1.
Commitment Level	Unstable

## **`_shutack_wait_bound`**

Description	Controls the maximum time, in milliseconds, to wait for a SHUTDOWN ACK after having sent a SHUTDOWN chunk.
Default	60,000
Range	0 to 300,000
Dynamic?	Yes
When to Change	Generally, you do not need to change this value. This parameter might be changed in accordance with “ <a href="#">_rto_max</a> ” on page 158.
Commitment Level	Unstable

## **`_maxburst`**

Description	Sets the limit on the number of segments to be sent in a burst.
Default	4
Range	2 to 8
Dynamic?	Yes

When to Change	You do not need to change this parameter. You might change it for testing purposes.
Commitment Level	Unstable

## **`_addip_enabled`**

Description	Enables or disables SCTP dynamic address reconfiguration.
Default	0 (disabled)
Range	0 (disabled) or 1 (enabled)
Dynamic?	Yes
When to Change	The parameter can be enabled if dynamic address reconfiguration is needed. Due to security implications, enable this parameter only for testing purposes.
Commitment Level	Unstable

## **`_prscpt_enabled`**

Description	Enables or disables the partial reliability extension (RFC 3758) to SCTP.
Default	1 (enabled)
Range	0 (disabled) or 1 (enabled)
Dynamic?	Yes
When to Change	Disable this parameter if partial reliability is not supported in your SCTP environment.
Commitment Level	Unstable

## **Per-Route Metrics**

Starting in the Solaris 8 release, you can use per-route metrics to associate some properties with IPv4 and IPv6 routing table entries.

For example, a system has two different network interfaces, a fast Ethernet interface and a gigabit Ethernet interface. The system default `_recv_hiwat` is 24,576 bytes. This default is sufficient for the fast Ethernet interface, but may not be sufficient for the gigabit Ethernet interface.



Instead of increasing the system's default for `_recv_hiwat`, you can associate a different default TCP receive window size to the gigabit Ethernet interface routing entry. By making this association, all TCP connections going through the route will have the increased receive window size.

For example, the following is in the routing table (`netstat -rn`), assuming IPv4:

```
192.123.123.0      192.123.123.4      U      1      4 hme0
192.123.124.0      192.123.124.4      U      1      4 ge0
default           192.123.123.1      UG     1      8
```

In this example, do the following:

```
# route change -net 192.123.124.0 -recvpipe x
```

Then, all connections going to the `192.123.124.0` network, which is on the `ge0` link, use the receive buffer size `x`, instead of the default 24567 receive window size.

If the destination is in the `a.b.c.d` network, and no specific routing entry exists for that network, you can add a prefix route to that network and change the metric. For example:

```
# route add -net a.b.c.d 192.123.123.1 -netmask w.x.y.z
# route change -net a.b.c.d -recvpipe y
```

Note that the prefix route's gateway is the default router. Then, all connections going to that network use the receive buffer size `y`. If you have more than one interface, use the `-ifp` argument to specify which interface to use. This way, you can control which interface to use for specific destinations. To verify the metric, use the `route(1M)` get command.



# Network Cache and Accelerator Tunable Parameters

---

This chapter describes some of the Network Cache and Accelerator (NCA) tunable parameters.

- “nca:nca\_conn\_hash\_size” on page 164
- “nca:nca\_conn\_req\_max\_q” on page 164
- “nca:nca\_conn\_req\_max\_q0” on page 164
- “nca:nca\_ppmax” on page 165
- “nca:nca\_vpmax” on page 165
- “sq\_max\_size” on page 166
- “ge:ge\_intr\_mode” on page 167

## Where to Find Tunable Parameters Information

Tunable Parameter	For Information
Solaris kernel tunables	Chapter 2, “Oracle Solaris Kernel Tunable Parameters”
NFS tunable parameters	Chapter 3, “NFS Tunable Parameters”
Internet Protocol Suite tunable parameters	Chapter 4, “Internet Protocol Suite Tunable Parameters”

## Tuning NCA Parameters

Setting these parameters is appropriate on a system that is a dedicated web server. These parameters allocate more memory for caching pages. You can set all of the tuning parameters described in this chapter in the `/etc/system` file.

For information on adding tunable parameters to the `/etc/system` file, see “[Tuning the Solaris Kernel](#)” on page 21.

## **nca:nca\_conn\_hash\_size**

Description	Controls the hash table size in the NCA module for all TCP connections, adjusted to the nearest prime number.
Default	383 hash table entries
Range	0 to 201,326,557
Dynamic?	No
When to Change	When the NCA's TCP hash table is too small to keep track of the incoming TCP connections. This situation causes many TCP connections to be grouped together in the same hashtable entry. This situation is indicated when NCA is receiving many TCP connections, and system performance decreases.
Commitment Level	Unstable

## **nca:nca\_conn\_req\_max\_q**

Description	Defines the maximum number of pending TCP connections for NCA to listen on.
Default	256 connections
Range	0 to 4,294,967,295
Dynamic?	No
When to Change	When NCA closes a connection immediately after it is established because it already has too many established TCP connections. If NCA is receiving many TCP connections and can handle a larger load, but is refusing any more connections, increase this parameter. Doing so allows NCA to handle more simultaneous TCP connections.
Commitment Level	Unstable

## **nca:nca\_conn\_req\_max\_q0**

Description	Defines the maximum number of incomplete (three-way handshake not yet finished) pending TCP connections for NCA to listen on.
Default	1024 connections
Range	0 to 4,294,967,295
Dynamic?	No

When to Change	When NCA refuses to accept any more TCP connections because it already has too many pending TCP connections. If NCA is receiving many TCP connections and can handle a larger load, but is refusing any more connections, increase this parameter. Doing so allows NCA to handle more simultaneous TCP connections.
Commitment Level	Unstable

## nca:nca\_ppmax

Description	Specifies the maximum amount of physical memory (in pages) used by NCA for caching the pages. This value should not be more than 75 percent of total memory.
Default	25 percent of physical memory
Range	1 percent to maximum amount of physical memory
Dynamic?	No
When to Change	When using NCA on a system with more than 512 MB of memory. If a system has a lot of physical memory that is not being used, increase this parameter. Then, NCA will efficiently use this memory to cache new objects. As a result, system performance will increase.  This parameter should be increased in conjunction with <code>nca_vpmax</code> , unless you have a system with more physical memory than virtual memory (a 32-bit kernel that has greater than 4 GB memory). Use <code>pagesize(1)</code> to determine your system's page size.
Commitment Level	Unstable

## nca:nca\_vpmax

Description	Specifies the maximum amount of virtual memory (in pages) used by NCA for caching pages. This value should not be more than 75 percent of the total memory.
Default	25 percent of virtual memory
Range	1 percent to maximum amount of virtual memory
Dynamic?	No
When to Change	When using NCA on a system with more than 512 MB of memory. If a system has a lot of virtual memory that is not being used, increase this

parameter. Then, NCA will efficiently use this memory to cache new objects. As a result, system performance will increase.

This parameter should be increased in conjunction with `nca_ppmax`. Set this parameter about the same value as `nca_vpmax`, unless you have a system with more physical memory than virtual memory.

Commitment Level    Unstable

## General System Tuning for the NCA

In addition to setting the NCA parameters, you can do some general system tuning to benefit NCA performance. If you are using gigabit Ethernet (ge driver), you should set the interface in interrupt mode for better results.

For example, a system with 4 GB of memory that is booted under 64-bit kernel should have the following parameters set in the `/etc/system` file. Use `pagesize` to determine your system's page size.

```
set sq_max_size=0
set ge:ge_intr_mode=1
set nca:nca_conn_hash_size=82500
set nca:nca_conn_req_max_q=100000
set nca:nca_conn_req_max_q0=100000
set nca:nca_ppmax=393216
set nca:nca_vpmax=393216
```

### sq\_max\_size

Description	Sets the depth of the syncq (number of messages) before a destination STREAMS queue generates a QFULL message.
Default	10000 messages
Range	0 (unlimited) to MAXINT
Dynamic?	No
When to Change	When NCA is running on a system with a lot of memory, increase this parameter to allow drivers to queue more packets of data. If a server is under heavy load, increase this parameter so that modules and drivers can process more data without dropping packets or getting backlogged.
Commitment Level	Unstable

## **ge:ge\_intr\_mode**

Description	Enables the ge driver to send packets directly to the upper communication layers rather than queue the packets
Default	0 (queue packets to upper layers)
Range	0 (enable) or 1 (disable)
Dynamic?	No
When to Change	When NCA is enabled, set this parameter to 1 so that the packet is delivered to NCA in interrupt mode for faster processing.
Commitment Level	Unstable





## System Facility Parameters

---

This chapter describes most of the parameters default values for various system facilities.

- “autofs” on page 170
- “cron” on page 170
- “devfsadm” on page 170
- “dhcagent” on page 170
- “fs” on page 170
- “ftp” on page 171
- “inetinit” on page 171
- “init” on page 171
- “ipsec” on page 171
- “kbd” on page 171
- “keyserv” on page 171
- “login” on page 171
- “mpathd” on page 172
- “nfs” on page 172
- “nfslogd” on page 172
- “nss” on page 172
- “passwd” on page 172
- “power” on page 172
- “su” on page 172
- “syslog” on page 172
- “tar” on page 173
- “utmpd” on page 173
- “yppasswdd” on page 173

# System Default Parameters

The functioning of various system facilities is governed by a set of values that are read by each facility on startup. The values stored in a file for each facility are located in the `/etc/default` directory. Not every system facility has a file located in this directory.

## **autofs**

This facility enables you to configure `autofs` parameters such as automatic timeout, displaying or logging status messages, browsing `autofs` mount points, and tracing. For details, see [autofs\(4\)](#).

## **cron**

This facility enables you to disable or enable `cron` logging.

## **devfsadm**

This file is not currently used.

## **dhcpgent**

Client usage of DHCP is provided by the `dhcpgent` daemon. When `ipadm` is used to create a DHCP address object, or when `ifconfig` identifies an interface that has been configured to receive its network configuration from DHCP, `dhcpgent` is started to manage an address on that interface.

For more information, see the `/etc/default/dhcpgent` information in the FILES section of [dhcpgent\(1M\)](#).

## **fs**

File system administrative commands have a generic and file system-specific portion. If the file system type is not explicitly specified with the `-F` option, a default is applied. The value is specified in this file. For more information, see the Description section of [default\\_fs\(4\)](#).

## ftp

This facility enables you to set the `ls` command behavior to the RFC 959 NLST command. The default `ls` behavior is the same as in the previous Solaris release.

For details, see [ftp\(4\)](#).

## inetinit

This facility enables you to configure TCP sequence numbers and to enable or disable support for 6to4 relay routers.

## init

For details, see the `/etc/default/init` information in the FILES section of [init\(1M\)](#).

All values in the file are placed in the environment of the shell that `init` invokes in response to a single user boot request. The `init` process also passes these values to any commands that it starts or restarts from the `/etc/inittab` file.

## ipsec

This facility enables you to configure parameters, such as IKE daemon debugging information and the `ikeadm` privilege level.

## kbd

For details, see the Extended Description section of [kbd\(1\)](#).

## keyserv

For details, see the `/etc/default/keyserv` information in the FILES section of [keyserv\(1M\)](#).

## login

For details, see the `/etc/default/login` information in the FILES section of [login\(1\)](#).

## **mpathd**

This facility enables you to set `in.mpathd` configuration parameters.

For details, see [in.mpathd\(1M\)](#).

## **nfs**

This facility enables you to set NFS daemon configuration parameters.

For details, see [nfs\(4\)](#).

## **nfslogd**

For details, see the Description section of [nfslogd\(1M\)](#).

## **nss**

This facility enables you to configure `initgroups(3C)` lookup parameters.

For details, see [nss\(4\)](#).

## **passwd**

For details, see the `/etc/default/passwd` information in the FILES section of [passwd\(1\)](#).

## **power**

For details, see the `/etc/default/power` information in the FILES section of [pmconfig\(1M\)](#).

## **su**

For details, see the `/etc/default/su` information in the FILES section of [su\(1M\)](#).

## **syslog**

For details, see the `/etc/default/syslogd` information in the FILES section of [syslogd\(1M\)](#).

## sys - suspend

For details, see the `/etc/default/sys-suspend` information in the FILES section of `sys-suspend(1M)`.

## tar

For a description of the `-f` function modifier, see `tar(1)`.

If the TAPE environment variable is not present and the value of one of the arguments is a number and `-f` is not specified, the number matching the `archiveN` string is looked up in the `/etc/default/tar` file. The value of the `archiveN` string is used as the output device with the blocking and size specifications from the file.

For example:

```
% tar -c 2 /tmp/*
```

This command writes the output to the device specified as `archive2` in the `/etc/default/tar` file.

## utmpd

The `utmpd` daemon monitors `/var/adm/utmpx` (and `/var/adm/utmp` in earlier Solaris versions) to ensure that `utmp` entries inserted by non-root processes by `pututxline(3C)` are cleaned up on process termination.

Two entries in `/etc/default/utmpd` are supported:

- `SCAN_PERIOD` – The number of seconds that `utmpd` sleeps between checks of `/proc` to see if monitored processes are still alive. The default is 300.
- `MAX_FDS` – The maximum number of processes that `utmpd` attempts to monitor. The default value is 4096 and should never need to be changed.

## yppasswdd

This facility enables you to configure whether a user can successfully set a login shell to a restricted shell when using the `passwd -r nis -e` command.

For details, see `rpc.yppasswdd(1M)`.



# Tunable Parameters Change History

---

This chapter describes the change history of specific tunable parameters. If a parameter is in this section, it has changed from a previous release. Parameters whose functionality has been removed are listed also.

- [“Kernel Parameters” on page 175](#)
- [“Parameters That Are Obsolete or Have Been Removed” on page 176](#)

## Kernel Parameters

### Process-Sizing Tunables

#### **ngroups\_max (Oracle Solaris 11 Express)**

This parameter was undocumented in previous Solaris releases. In this Solaris release, the default maximum has been increased to 1024 groups. For more information, see [“ngroups\\_max” on page 42](#).

### General Driver Parameter

#### **ddi\_msix\_alloc\_limit (Solaris 10 Release and Oracle Solaris 11 Express Release)**

This parameter is new starting in the Solaris 10 10/09 release and the Oracle Solaris 11 Express release. For more information, see [“ddi\\_msix\\_alloc\\_limit” on page 59](#).

## Network Driver Parameters

### **igb Parameters (Oracle Solaris 11 Express Release)**

The `igb` network driver parameters are provided in the Oracle Solaris 11 Express release. For more information, see [“igb Parameters” on page 60](#).

### **ixgbe Parameters (Oracle Solaris 11 Express Release)**

The `ixgbe` network driver parameters are provided in the Oracle Solaris 11 Express release. For more information, see [“ixgbe Parameters” on page 61](#).

## General Kernel and Memory Parameters

### **zfs\_arc\_min (Oracle Solaris 11 Express)**

This parameter description is newly documented in the Solaris 10 10/09 release. For more information, see [“zfs\\_arc\\_min” on page 29](#).

### **zfs\_arc\_max (Oracle Solaris 11 Express)**

This parameter description is newly documented in the Solaris 10 10/09 release. For more information, see [“zfs\\_arc\\_max” on page 29](#).

## fsflush and Related Parameters

### **dopageflush (All Solaris Releases)**

The description was clarified by including that number of *physical* memory pages are examined.

## Parameters That Are Obsolete or Have Been Removed

The following section describes parameters that are obsolete or have been removed from more recent Solaris releases.

## TCP/IP Module Parameters

### **ip\_multidata\_outbound (Oracle Solaris 11 Express)**

This parameter is obsolete in the Oracle Solaris 11 Express release.



## **tcp\_mdt\_max\_pbufs (Oracle Solaris 11 Express)**

This parameter is obsolete in the Oracle Solaris 11 Express release.



## Revision History for This Manual

---

This section describes the revision history for this manual.

- [“Current Version: Oracle Solaris 11 Express Release” on page 179](#)

### Current Version: Oracle Solaris 11 Express Release

The current version of this manual applies to the Oracle Solaris 11 Express release.

### New or Changed Parameters in the Oracle Solaris Release

The following sections describe new, changed, or obsolete kernel tunables.

- Oracle Solaris 11 Express: The `ip_multidata_outbound` parameter and the `tcp_mdt_max_pbufs` parameter for devices that support multidata transport (MDT) are obsolete in this release.
- Oracle Solaris 11 Express: This release includes the `ngroups_max` parameter description. For more information, see [“ngroups\\_max” on page 42](#).
- Oracle Solaris 11 Express: This release includes the `zfs_arc_min` and `zfs_arc_max` parameter descriptions. For more information, see [“zfs\\_arc\\_min” on page 29](#) and [“zfs\\_arc\\_max” on page 29](#).
- Oracle Solaris 11 Express: This release includes several `igb` and `ixgbe` network driver parameters. For more information, see [“igb Parameters” on page 60](#) and [“ixgbe Parameters” on page 61](#).
- Oracle Solaris 11 Express: This release includes the `ddi_msix_alloc_limit` parameter that can be used to increase the number of MSI-X interrupts that a device instance can allocate. For more information, see [“ddi\\_msix\\_alloc\\_limit” on page 59](#).
- Oracle Solaris 11 Express: This release includes corrected range information for the `tcp_local_dack_interval` parameter. For more information, see [“\\_local\\_dack\\_interval” on page 137](#).

- Oracle Solaris 11 Express: This release includes the `kmem_stackinfo` parameter, which can be enabled to monitor kernel thread stack usage. For more information, see “[kmem\\_stackinfo](#)” on page 57.
- Oracle Solaris 11 Express: For information about tuning ZFS file systems, see the following site:  
[http://www.solarisinternals.com/wiki/index.php/ZFS\\_Evil\\_Tuning\\_Guide](http://www.solarisinternals.com/wiki/index.php/ZFS_Evil_Tuning_Guide)
- Oracle Solaris 11 Express: Memory locality group parameters are provided in this release. For more information about these parameters, see “[Locality Group Parameters](#)” on page 90.
- Oracle Solaris 11 Express: Parameter information was updated to include sun4v systems. For more information, see the following references:
  - “[maxphys](#)” on page 65
  - “[tmpfs:tmpfs\\_maxkmem](#)” on page 78
  - “[sun4u or sun4v Specific Parameters](#)” on page 86

# Index

---

## A

\_addip\_enabled, 160  
\_addrs\_per\_if, 135  
autofs, 170  
autoup, 36

## B

bufhwm, 71  
bufhwm\_pct, 71

## C

\_conn\_req\_max\_q, 144  
\_conn\_req\_max\_q0, 144  
\_conn\_req\_min, 145  
consistent\_coloring, 86  
\_cookie\_life, 158  
cron, 170  
\_cwnd\_max, 140, 154

## D

ddi\_msix\_alloc\_limit parameter, 59  
default\_stksize, 30  
default\_tsb\_size, 88  
\_deferred\_ack\_interval, 136, 156  
\_deferred\_acks\_max, 137  
desfree, 45  
dhcpagent, 170

dnlc\_dir\_enable, 69  
dnlc\_dir\_max\_size, 70  
dnlc\_dir\_min\_size, 69  
doiflush, 37  
dopageflush, 37, 176

## E

ecn, 143  
enable\_tsb\_rss\_sizing, 89

## F

fastscan, 50  
freebehind, 76  
fs, 170  
fsflush, 34  
ftp, 171

## G

ge\_intr\_mode, 167

## H

handspreadpages, 51  
\_heartbeat\_interval, 155  
hires\_tick, 85

**I**

\_icmp\_err\_burst, 135  
\_icmp\_err\_interval, 135  
\_ignore\_path\_mtu, 156  
inetinit, 171  
init, 171  
\_initial\_mtu, 155  
\_initial\_out\_streams, 159  
\_initial\_ssthresh, 156  
intr\_force, 60  
intr\_throttling, 62  
\_ip\_abort\_interval, 148  
ip\_squeue\_fanout, 136  
ip\_squeue\_worker\_wait, 147  
ipcl\_conn\_hash\_size, 146  
ipsec, 171  
\_ipv4\_ttl, 154  
\_ipv6\_hoplimit, 157

**K**

kbd, 171  
\_keepalive\_interval, 148  
keyserv, 171  
kmem\_flags, 55  
kmem\_stackinfo, 57

**L**

lgrp\_mem\_pset\_aware, 92  
\_local\_dack\_interval, 137  
\_local\_dacks\_max, 138  
logevent\_max\_q\_sz, 32  
login, 171  
lotsfree, 44  
lpg\_alloc\_prefer, 90  
lpg\_mem\_default\_policy, 91  
lwp\_default\_stksize, 31

**M**

\_max\_buf, 140, 157

\_max\_in\_streams, 158  
\_max\_init\_retr, 153  
max\_nprocs, 41  
maxpgio, 53  
maxphys, 65  
maxpid, 40  
maxuprc, 41  
maxusers, 38  
min\_percent\_cpu, 51  
minfree, 46  
moddebug, 58  
mpathd, 172  
mr\_enable, 60

**N**

nca\_conn\_hash\_size, 164  
nca\_conn\_req\_max\_q, 164  
nca\_conn\_req\_max\_q0, 164  
nca\_ppmax, 165  
nca\_vpmax, 165  
ncsize, 67  
ndd, 134  
ndquot, 73  
\_new\_secret\_interval, 155  
nfs\_max\_threads, 105  
nfs:nacache, 119  
nfs:nfs\_allow\_preepoch\_time, 97  
nfs:nfs\_async\_clusters, 116  
nfs:nfs\_async\_timeout, 118  
nfs:nfs\_bsize, 114  
nfs:nfs\_cots\_timeo, 98  
nfs:nfs\_disable\_rddir\_cache, 113  
nfs:nfs\_do\_symlink\_cache, 100  
nfs:nfs\_dynamic, 102  
nfs:nfs\_lookup\_neg\_cache, 103  
nfs:nfs\_nra, 107  
nfs:nfs\_shrinkreaddir, 110  
nfs:nfs\_write\_error\_interval, 112  
nfs:nfs\_write\_error\_to\_cons\_only, 112  
nfs:nfs3\_async\_clusters, 117  
nfs:nfs3\_bsize, 114  
nfs:nfs3\_cots\_timeo, 98  
nfs:nfs3\_do\_symlink\_cache, 100

nfs:nfs3\_dynamic, 102  
 nfs:nfs3\_jukebox\_delay, 120  
 nfs:nfs3\_lookup\_neg\_cache, 103  
 nfs:nfs3\_max\_threads, 106  
 nfs:nfs3\_max\_transfer\_size, 120  
 nfs:nfs3\_max\_transfer\_size\_clts, 122  
 nfs:nfs3\_max\_transfer\_size\_cots, 123  
 nfs:nfs3\_nra, 108  
 nfs:nfs3\_pathconf\_disable\_cache, 96  
 nfs:nfs3\_shrinkreaddir, 111  
 nfs:nfs4\_async\_clusters, 118  
 nfs:nfs4\_bsize, 115  
 nfs:nfs4\_cots\_timeo, 99  
 nfs:nfs4\_do\_symlink\_cache, 101  
 nfs:nfs4\_lookup\_neg\_cache, 104  
 nfs:nfs4\_max\_threads, 107  
 nfs:nfs4\_max\_transfer\_size, 121  
 nfs:nfs4\_nra, 109  
 nfs:nfs4\_pathconf\_disable\_cache, 96  
 nfs:nrnode, 109  
 nfslogd, 172  
 nfssrv:exi\_cache\_time, 126  
 nfssrv:nfs\_portmon, 123  
 nfssrv:nfsauth\_ch\_cache\_max, 125  
 nfssrv:rfs\_write\_async, 124  
 ngroups\_max, 42, 175  
 noexec\_user\_stack, 34  
 nss, 172  
 nstrpush, 82

## P

pageout\_reserve, 47  
 pages\_before\_pager, 52  
 pages\_pp\_maximum, 48  
 passwd, 172  
 physmem, 28  
 pidmax, 40  
 \_policy\_mask, 152  
 power, 172  
 \_pp\_max\_retr, 154  
 \_prsectp\_enabled, 160  
 pt\_cnt, 80  
 pt\_max\_pty, 81

pt\_pctofmem, 80

## R

rechoose\_interval, 85  
 \_recv\_hiwat\_minmss, 151  
 recv\_maxbuf, 139, 151  
 reserved\_procs, 39  
 \_respond\_to\_echo\_broadcast, 135  
 \_rev\_src\_routes, 142  
 \_rexmit\_interval\_extra, 150  
 \_rexmit\_interval\_initial, 149  
 \_rexmit\_interval\_max, 149  
 \_rexmit\_interval\_min, 149  
 rlim\_fd\_cur, 66  
 rlim\_fd\_max, 65  
 rpcmod:clnt\_idle\_timeout, 127  
 rpcmod:clnt\_max\_conns, 126  
 rpcmod:cotsmaxdupreqs, 130  
 rpcmod:maxdupreqs, 129  
 rpcmod:svc\_default\_stksize, 128  
 rpcmod:svc\_idle\_timeout, 127  
 \_rst\_sent\_rate, 146  
 \_rst\_sent\_rate\_enabled, 145  
 rstchown, 68  
 \_rto\_max, 158  
 \_rto\_min, 157  
 rx\_copy\_threshold, 64  
 rx\_limit\_per\_intr, 62  
 rx\_queue\_number, 61  
 rx\_ring\_size, 63

## S

sack, 141  
 sctp\_maxburst, 159  
 segmap\_percent, 70  
 segspt\_minfree, 84  
 send\_maxbuf, 139, 151  
 \_shutack\_wait\_bound, 159  
 \_slow\_start\_after\_idle, 141  
 \_slow\_start\_initial, 141  
 slowscan, 50

smallfile, 77  
sq\_max\_size, 166  
strmsgsz, 82,83  
su, 172  
sun4u, 86  
sun4v, 86  
swapfs\_minfree, 54  
swapfs\_reserve, 54  
sys-suspend, 173  
syslog, 172

## T

tar, 173  
throttlefree, 47  
\_time\_wait\_interval, 142  
timer\_max, 86  
tmpfs\_maxkmem, 78  
tmpfs\_minfree, 78  
tsb\_alloc\_hiwater, 87  
tsb\_rss\_size, 89  
\_tstamp\_always, 139  
\_tstamp\_if\_wscale, 150  
tune\_t\_fsflushr, 35  
tune\_t\_minarmem, 49  
tx\_copy\_threshold, 64  
tx\_queue\_number, 61  
tx\_ring\_size, 63

## U

ufs\_HW, 75  
ufs\_LW, 75  
ufs\_ninode, 73  
ufs:ufs\_WRITES, 75  
utmpd, 173

## W

\_wscale\_always, 138

## Y

yppasswdd, 173

## Z

zfs\_arc\_max, 29,176  
zfs\_arc\_min, 29,176