Oracle Linux 9 Managing Core System Configuration





Oracle Linux 9 Managing Core System Configuration,

F56701-06

Copyright $\ensuremath{\texttt{@}}$ 2022, 2023, Oracle and/or its affiliates.

Contents

Preface

Conventions	V
Documentation Accessibility	V
Access to Oracle Support for Accessibility	V
Diversity and Inclusion	V
Managing Kernels and System Boot	
About the Boot Process	1-1
About UEFI-Based Booting	1-1
About BIOS-Based Booting	1-2
About the GRUB 2 Bootloader	1-2
About Linux Kernels	1-3
Managing Kernels in GRUB 2 Using grubby	1-4
Kernel Boot Parameters	1-5
Modifying Kernel Boot Parameters Before Booting	1-7
	1-7 1-7
Modifying GRUB 2 Default Kernel Boot Parameters Managing System Services With systemd About the systemd Service Manager	2-2
Modifying GRUB 2 Default Kernel Boot Parameters Managing System Services With systemd About the systemd Service Manager systemd Units	2-3 2-3
Modifying GRUB 2 Default Kernel Boot Parameters Managing System Services With systemd About the systemd Service Manager systemd Units About System-State Targets	2-2 2-2 2-2
Modifying GRUB 2 Default Kernel Boot Parameters Managing System Services With systemd About the systemd Service Manager systemd Units About System-State Targets Displaying Default and Active System-State Targets	2-3 2-3 2-3 2-3
Managing System Services With systemd About the systemd Service Manager systemd Units About System-State Targets Displaying Default and Active System-State Targets Changing Default and Active System-State Targets	2-1 2-2 2-3 2-4
Managing System Services With systemd About the systemd Service Manager systemd Units About System-State Targets Displaying Default and Active System-State Targets Changing Default and Active System-State Targets Shutting Down, Suspending, and Rebooting the System	2-1 2-3 2-3 2-4 2-5
Managing System Services With systemd About the systemd Service Manager systemd Units About System-State Targets Displaying Default and Active System-State Targets Changing Default and Active System-State Targets Shutting Down, Suspending, and Rebooting the System Managing Services	2-2 2-3 2-4 2-4 2-5 2-5 2-5
Managing System Services With systemd About the systemd Service Manager systemd Units About System-State Targets Displaying Default and Active System-State Targets Changing Default and Active System-State Targets Shutting Down, Suspending, and Rebooting the System Managing Services Starting and Stopping Services	2-2-2-3 2-4 2-5 2-5 2-5 2-5 2-5
Managing System Services With systemd About the systemd Service Manager systemd Units About System-State Targets Displaying Default and Active System-State Targets Changing Default and Active System-State Targets Shutting Down, Suspending, and Rebooting the System Managing Services Starting and Stopping Services Enabling and Disabling Services	2-3 2-3 2-3 2-4 2-4 2-5 2-5 2-6
Managing System Services With systemd About the systemd Service Manager systemd Units About System-State Targets Displaying Default and Active System-State Targets Changing Default and Active System-State Targets Shutting Down, Suspending, and Rebooting the System Managing Services Starting and Stopping Services Enabling and Disabling Services Displaying the Status of Services	2-5 2-5 2-5 2-5 2-5 2-5 2-5 2-5 2-5 2-5
Managing System Services With systemd About the systemd Service Manager systemd Units About System-State Targets Displaying Default and Active System-State Targets Changing Default and Active System-State Targets Shutting Down, Suspending, and Rebooting the System Managing Services Starting and Stopping Services Enabling and Disabling Services	2-1 2-2 2-3



	About Service Unit Files	2-10
	Configurable Options in Service Unit Files	2-11
	Creating a User-Based systemd Service	2-13
3	Configuring System Settings	
	About the /etc/sysconfig Files	3-1
	About the /proc Virtual File System	3-2
	Virtual Files and Directories Under /proc	3-3
	Modifying Kernel Parameters	3-8
	Parameters That Control System Performance	3-9
	Parameters That Control Kernel Panics	3-10
	About the /sys Virtual File System	3-12
	Virtual Directories Under the /sys Directory	3-12
	Configuring System Date and Time Settings	3-13
4	Managing System Devices	
	About Device Files	4-1
	About the Udev Device Manager	4-3
	About Udev Rules	4-3
	Querying Udev and Sysfs	4-7
	Modifying Udev Rules	4-10
5	Managing Kernel Modules	
	About Kernel Modules	5-1
	Listing Information About Loaded Modules	5-1
	Loading and Unloading Modules	5-3
	About Module Parameters	5-4
	Specifying Modules To Be Loaded at Boot Time	5-5
	Preventing Modules From Loading at Boot Time	5-5
	About Weak Update Modules	5-6
6	Configuring Huge Pages	
	Available Huge Page Features	6-1
	HugeTLB Pages	6-1
	Transparent HugePages	6-2
	Configuring HugeTLB Pages	6-2
	Kernel Boot Parameters for HugeTLB Pages	6-2
	File-Based Configuration Parameters for HugeTLB Pages	6-3



Configuring HugeTLB Pages at Boot Time	6-9
Requesting HugeTLB Pages by Using Kernel Parameters at Boot Time	6-9
Requesting HugeTLB Pages Using NUMA Node-Specific Parameters Early in the Boot Process	6-9
Configuring HugeTLB at Runtime	6-10
Configuring HugeTLB Pages for a Specific NUMA Node at Runtime	6-10
Configuring Transparent HugePages	6-11
Parameters Used to Configure Transparent HugePages	6-11
Configuring Transparent HugePages at Runtime	6-14
Retrieving the Current Status of Transparent HugePages	6-14
Changing the Current Status of Transparent HugePages	6-14
Changing the defrag Setting of Transparent HugePages	6-15
Managing Resources	
	7-1
About Control Groups	
About Control Groups About Kernel Resource Controllers	7-1
About Control Groups About Kernel Resource Controllers About the Control Group File System	7-1 7-2
About Control Groups About Kernel Resource Controllers About the Control Group File System About Control Groups and systemd	7-1 7-2 7-2
About Control Groups About Kernel Resource Controllers About the Control Group File System About Control Groups and systemd About Resource Distribution Models	7-1 7-1 7-2 7-2 7-4 7-4
About Control Groups About Kernel Resource Controllers About the Control Group File System About Control Groups and systemd About Resource Distribution Models	7-1 7-2 7-2 7-4
About Control Groups About Kernel Resource Controllers About the Control Group File System About Control Groups and systemd About Resource Distribution Models Using cgroups v2 to Manage Resources for Applications	7-1 7-2 7-2 7-4 7-4 7-5
About Control Groups About Kernel Resource Controllers About the Control Group File System About Control Groups and systemd About Resource Distribution Models Using cgroups v2 to Manage Resources for Applications Enabling cgroups v2	7-1 7-2 7-2 7-4 7-4
Preparing the Control Group for Distribution of CPU Time	7-1 7-2 7-2 7-4 7-4 7-5
About Control Groups About Kernel Resource Controllers About the Control Group File System About Control Groups and systemd About Resource Distribution Models Using cgroups v2 to Manage Resources for Applications Enabling cgroups v2 Preparing the Control Group for Distribution of CPU Time Setting CPU Bandwidth to Regulate Distribution of CPU Time	7-1 7-2 7-2 7-4 7-5 7-5



Preface

Oracle Linux 9: Managing Core System Configuration provides information about configuring Oracle Linux 9 systems, including the boot loader configuration and processes, system devices, services and settings, as well as kernel parameters.

Conventions

The following text conventions are used in this document:

Convention	Meaning
boldface	Boldface type indicates graphical user interface elements associated with an action, or terms defined in text or the glossary.
italic	Italic type indicates book titles, emphasis, or placeholder variables for which you supply particular values.
monospace	Monospace type indicates commands within a paragraph, URLs, code in examples, text that appears on the screen, or text that you enter.

Documentation Accessibility

For information about Oracle's commitment to accessibility, visit the Oracle Accessibility Program website at https://www.oracle.com/corporate/accessibility/.

For information about the accessibility of the Oracle Help Center, see the Oracle Accessibility Conformance Report at https://www.oracle.com/corporate/accessibility/templates/t2-11535.html.

Access to Oracle Support for Accessibility

Oracle customers that have purchased support have access to electronic support through My Oracle Support. For information, visit https://www.oracle.com/corporate/accessibility/learning-support.html#support-tab.

Diversity and Inclusion

Oracle is fully committed to diversity and inclusion. Oracle respects and values having a diverse workforce that increases thought leadership and innovation. As part of our initiative to build a more inclusive culture that positively impacts our employees, customers, and partners, we are working to remove insensitive terms from our

products and documentation. We are also mindful of the necessity to maintain compatibility with our customers' existing technologies and the need to ensure continuity of service as Oracle's offerings and industry standards evolve. Because of these technical constraints, our effort to remove insensitive terms is ongoing and will take time and external cooperation.



1

Managing Kernels and System Boot

This chapter describes the Oracle Linux boot process and how to configure and use the GRand Unified Bootloader (GRUB) version 2 and boot-related kernel parameters.



Tip:

See Manage the Boot Kernel for Oracle Linux for a hands-on tutorial and video demonstrations on configuring the boot kernel in Oracle Linux.

About the Boot Process

Understanding the Oracle Linux boot process can help you troubleshoot problems when booting a system. The boot process involves several files, and errors in these files are the usual cause of boot problems. Boot processes and configuration differ depending on whether your hardware uses UEFI firmware or legacy BIOS to handle system boot.

About UEFI-Based Booting

On a UEFI-based system running the Oracle Linux release, the system boot process uses the following sequence:

- The system's UEFI firmware performs a power-on self-test (POST) and then locates and initializes peripheral devices and the hard disk.
- 2. UEFI searches for a GPT partition with a specific globally unique identifier (GUID) that identifies it as the EFI System Partition (ESP), which contains EFI applications such as boot loaders. In case of the presence of multiple boot devices, the UEFI boot manager determines the appropriate ESP to use, based on the order that is defined in the boot manager. With the efibootmgr tool, you can define a different order, if you do not want to use the default definition.
- The UEFI boot manager checks to determine whether Secure Boot is enabled. If Secure Boot is not enabled, the boot manager runs the GRUB 2 bootloader on the ESP.
 - Otherwise, the boot manager requests a certificate from the boot loader and validates this against keys stored in the UEFI Secure Boot key database. To handle the certificate validation process, the environment is configured to perform a 2-stage boot process and the <code>shim.efi</code> application that is responsible for certification is loaded first before loading the GRUB 2 bootloader. If the certificate is valid, the boot loader runs and, in turn, validates the kernel that it is configured to load. See Oracle Linux: Working With UEFI Secure Boot for more information on Secure Boot.
- 4. The boot loader loads the vmlinuz kernel image file into memory and extracts the contents of the initramfs image file into a temporary, memory-based file system (tmpfs).
- 5. The kernel loads the driver modules from the initramfs file system that are needed to access the root file system.

- 6. The kernel starts the systemd process with a process ID of 1 (PID 1). See About the systemd Service Manager.
- 7. systemd runs any additional processes defined for it.



Specify any other actions to be processed during the boot process by defining your own systemd unit. This method is the preferred approach than using the /etc/rc.local file.

About BIOS-Based Booting

On a BIOS-based system running the Oracle Linux release, the boot process is as follows:

- 1. The system's BIOS performs a power-on self-test (POST), and then locates and initializes any peripheral devices and the hard disk.
- 2. The BIOS reads the Master Boot Record (MBR) into memory from the boot device. The MBR stores information about the organization of partitions on that device, the partition table, and the boot signature which is used for error detection. Additionally, the MBR includes the pointer to the boot loader program (GRUB 2). The boot program itself can be on the same device or on another device.
- The boot loader loads the vmlinuz kernel image file into memory and extracts the contents of the initramfs image file into a temporary, memory-based file system (tmpfs).
- 4. The kernel loads the driver modules from the initramfs file system that are needed to access the root file system.
- 5. The kernel starts the systemd process with a process ID of 1 (PID 1). See About the systemd Service Manager for more information.
- 6. systemd runs any additional processes defined for it.

Note:

Specify any other actions to be processed during the boot process by defining your own systemd unit. This method is the preferred approach than using the /etc/rc.local file.

About the GRUB 2 Bootloader

GRUB 2 can load many operating systems in addition to Oracle Linux and it can chain-load proprietary operating systems. GRUB 2 understands the formats of file systems and kernel executables, which allows it to load an arbitrary operating system without needing to know the exact location of the kernel on the boot device. GRUB 2 requires only the file name and drive partitions to load a kernel. You can configure this information by using the GRUB 2 menu or by entering it on the command line.



GRUB 2 behavior is based on configuration files. On BIOS-based systems, the configuration file is /boot/grub2/grub.cfg. On UEFI-based systems, the configuration file is /boot/efi/EFI/redhat/grub.cfg. Each kernel version's boot paramaters are stored in independent configuration files in /boot/loader/entries. Each kernel configuration is stored with the file name $machine\ id-kernel\ version.el8.arch.conf.$



Do not edit the GRUB 2 configuration file directly.

The grub2-mkconfig command generates the configuration file using the template scripts in /etc/grub.d and menu-configuration settings taken from the configuration file, /etc/default/grub.

The default menu entry is determined by the value of the <code>GRUB_DEFAULT</code> parameter in /etc/default/grub. Setting <code>GRUB_DEFAULT</code> to saved allows you to use the <code>grub2-set-default</code> and <code>grub2-reboot</code> commands to specify the default entry. The command <code>grub2-set-default</code> sets the default entry for all subsequent reboots and <code>grub2-reboot</code> sets the default entry for the next reboot only.

If you specify a numeric value as the value of <code>GRUB_DEFAULT</code> or as an argument to either <code>grub2-reboot</code> or <code>grub2-set-default</code>, GRUB 2 counts the menu entries in the configuration file starting at 0 for the first entry.

For more information about using, configuring, and customizing GRUB 2, see the GNU GRUB Manual, which is also installed as /usr/share/doc/grub2-tools-2.00/grub.html.

About Linux Kernels

The Linux Foundation provides a hub for open source developers to code, manage, and scale a variety of open technology projects. It also manages the Linux Kernel Organization that exists to distribute various versions of the Linux kernel which is at the core of all Linux distributions, including those used by Oracle Linux. The Linux kernel manages the interactions between the computer hardware and user space applications that run on Oracle Linux.

You must install and run one of these Linux kernels with Oracle Linux:

- Unbreakable Enterprise Kernel (UEK): UEK is based on a stable kernel branch from the Linux Foundation, with customer-driven additions, and multiple UEKs can exist for a given Oracle Linux release. Its focus is performance, stability, and minimal backports by tracking the mainline source code provided by the Linux Kernel Organization, as closely as is practical. UEK is well-tested and used to run Oracle's Engineered Systems, Oracle Cloud Infrastructure (OCI), and large enterprise deployments for Oracle customers. UEK includes some packages or package versions that are not available in RHCK. Some examples are btrfs-tools, rds, and rdma related packages, and some kernel tuning tools.
- Red Hat Compatible Kernel (RHCK): RHCK is identical to the Linux kernel that is
 distributed in a corresponding Red Hat Enterprise Linux (RHEL) release. You can use
 RHCK to ensure full compatibility with applications that run on Red Hat Enterprise Linux.



Important:

Linux kernels are critical for running applications in the Oracle Linux user space. Therefore, you must keep the kernel up to date with the latest bug fixes, enhancements, and security updates provided by Oracle. To do so, implement a continuous update and upgrade strategy. See Oracle Linux: Ksplice User's Guide for information on how to keep your kernel up to date without any requirement to reboot the system. See Oracle Linux: Managing Software on Oracle Linux for general information about keeping software on your system up to date.

See Unbreakable Enterprise Kernel documentation for more information about UEK.

Managing Kernels in GRUB 2 Using grubby

You can use the grubby command to view and manage your kernels.

Use the following command to display all of the kernels that are installed and configured on your system:

```
sudo grubby --info=ALL
```

To configure a specific kernel as the default boot kernel, run:

```
sudo grubby --set-default /boot/vmlinuz-4.18.0-80.el8.x86_64
```

You can also use the grubby command to update a kernel configuration entry to add or remove kernel boot arguments, for example:

```
sudo grubby --remove-args="rhgb quiet" --
args=rd_LUKS_UUID=luks-39fec799-6a6c-4ac1-ac7c-1d68f2e6b1a4 \
--update-kernel /boot/vmlinuz-4.18.0-80.el8.x86 64
```

For more information about the grubby command, see the grubby (8) manual page, which you can access by running the man grubby command.

For a hands-on tutorial that demonstrates the use of grubby to manage kernels, see Manage the Boot Kernel for Oracle Linux.



Important:

Security scanners on your system might report CVEs for any kernel on the system that is not exclusively used as the running or default kernel. As a good practice and to avoid unnecessary noise and false positives that are being reported by your scanner, remove unused kernels after you switch kernels.

For example, if you switch to RHCK from UEK, follow these steps to ensure a proper transition:

- Set the default kernel to RHCK with the appropriate grubby command.
 See the preceding examples for the correct command syntax.
- 2. Reboot the system to ensure that you are now running RHCK.
- 3. Remove the UEK kernel.

```
sudo dnf remove kernel-uek
```

4. Disable the UEK repositories by running the following command:

```
for uek_repo in (dnf repolist enabled|grep UEK|awk '{print $1}'); do sudo dnf config-manager disable <math>uek_repo; done
```

5. Downgrade kernel plumbing packages and remove orphan packages.

```
sudo dnf downgrade $(sudo package-cleanup --orphans)
```

Likewise, if you choose to use UEK as your standard kernel, consider removing RHCK to similarly avoid false positive warnings from your security scanner. For instructions to remove RHCK, see Remove the Red Hat Compatible Kernel With the kernel-transition Package.

Kernel Boot Parameters

The following table describes some commonly used kernel boot parameters.

Option	Description
0, 1, 2, 3, 4, 5, or 6, or	Specifies the nearest systemd-equivalent
systemd.unit=runlevelN.target	system-state target to match a legacy SysV run level. <i>N</i> can take an integer value between 0 and 6.
	Systemd maps system-state targets to mimic the legacy SysV init system. For a description of system-state targets, see About System-State Targets.
<pre>1, s, S, single, or systemd.unit=rescue.target</pre>	Specifies the rescue shell. The system boots to single-user mode prompts for the root password.
3 or systemd.unit=multi-user.target	Specifies the systemd target for multi-user, non-graphical login.
5 or systemd.unit=graphical.target	Specifies the systemd target for multi-user, graphical login.



Option	Description
-b, emergency, or systemd.unit=emergency.target	Specifies emergency mode. The system boots to single-user mode and prompts for the root password. Fewer services are started than when in rescue mode.
KEYBOARDTYPE=kbtype	Specifies the keyboard type, which is written to /etc/sysconfig/keyboard in the initramfs.
KEYTABLE=kbtype	Specifies the keyboard layout, which is written to /etc/sysconfig/keyboard in the initramfs.
LANG=language_territory.codeset	Specifies the system language and code set, which is written to /etc/sysconfig/i18n in the initramfs.
max_loop=N	Specifies the number of loop devices ($/\text{dev}/1\text{oop*}$) that are available for accessing files as block devices. The default and maximum values of N are 8 and 255.
nouptrack	Disables Ksplice Uptrack updates from being applied to the kernel.
quiet	Reduces debugging output.
rd_LUKS_UUID=UUID	Activates an encrypted Linux Unified Key Setup (LUKS) partition with the specified UUID.
rd_LVM_VG=vg/lv_vol	Specifies an LVM volume group and volume to be activated.
rd_NO_LUKS	Disables detection of an encrypted LUKS partition.
rhgb	Specifies that the Red Hat graphical boot display should be used to indicate the progress of booting.
rn_NO_DM	Disables Device-Mapper (DM) RAID detection.
rn_NO_MD	Disables Multiple Device (MD) RAID detection.
ro root=/dev/mapper/vg-lv_root	Specifies that the root file system is to be mounted read only, and specifies the root file system by the device path of its LVM volume (where vg is the name of the volume group).
rw root=UUID= <i>UUID</i>	Specifies that the root (/) file system is to be mounted read-writable at boot time, and specifies the root partition by its UUID.
selinux=0	Disables SELinux.
SYSFONT=font	Specifies the console font, which is written to /etc/sysconfig/il8n in the initramfs.

The kernel boot parameters that were last used to boot a system are recorded in / proc/cmdline, for example:

sudo cat /proc/cmdline

BOOT_IMAGE=(hd0,msdos1)/vmlinuz-4.18.0-80.el8.x86_64 root=/dev/mapper/ol-root ro



```
crashkernel=auto \
resume=/dev/mapper/ol-swap rd.lvm.lv=ol/root rd.lvm.lv=ol/swap rhgb quiet
```

For more information, see the kernel-command-line(7) manual page.

Modifying Kernel Boot Parameters Before Booting

To modify boot parameters before booting a kernel, follow these steps:

- When the GRUB boot menu appears at the beginning of the boot process, use the arrow keys to highlight the required kernel and press the space bar.
- 2. Press E to edit the boot configuration for the kernel.
- 3. Use the arrow keys to bring the cursor to the end of the line that starts with linux, which is the boot configuration line for the kernel.
- 4. Modify the boot parameters.

You can add parameters such as systemd.target=runlevel1.target, which instructs the system to boot into the rescue shell.

Press Ctrl+X to boot the system.

Modifying GRUB 2 Default Kernel Boot Parameters

To modify the boot parameters for the GRUB 2 configuration so that these parameters are applied by default at every reboot, follow these steps:

1. Edit /etc/default/grub and modify the parameters in the GRUB_CMDLINE_LINUX definition, for example:

```
GRUB_CMDLINE_LINUX="vconsole.font=latarcyrheb-sun16 vconsole.keymap=uk
crashkernel=auto rd.lvm.lv=ol/swap rd.lvm.lv=ol/root biosdevname=0
rhgb quiet systemd.unit=runlevel3.target"
```

This example adds the parameter systemd.unit=runlevel3.target so that the system boots into multi-user, non-graphical mode by default.

2. Rebuild /boot/grub2/grub.cfg:

```
sudo grub2-mkconfig -o /boot/grub2/grub.cfg
```

The change takes effect for subsequent system reboots of all configured kernels.

Note:

For systems that boot with UEFI, the <code>grub.cfg</code> file is located in the <code>/boot/efi/EFI/redhat</code> directory because the boot configuration is stored on a dedicated FAT32-formatted partition.

After the system has successfully booted, the EFI folder on that partition is mounted inside the /boot/efi directory on the root file system for Oracle Linux.



2

Managing System Services With systemd

The systemd daemon is the system initialization and service manager in Oracle Linux. This chapter describes how to use systemd to manage system processes, services and systemd targets.



Tip

See Use systemd on Oracle Linux for a hands-on tutorial and video demonstrations on working with systemd in Oracle Linux.

About the systemd Service Manager

The systemd daemon is the first process that starts after a system boots and is the final process that is running when the system shuts down. systemd controls the final stages of booting and prepares the system for use. It also speeds up booting by loading services concurrently.

systemd reads its configuration from files in the /etc/systemd directory. For example, the /etc/systemd/system.conf file controls how systemd handles system initialization.

To determine which services to start during the boot process, <code>systemd</code> reads the symbolic link <code>/etc/systemd/system/default.target</code>. The following example shows the value of <code>/etc/systemd/system/default.target</code> on a system configured to boot to a multi-user mode without a graphical user interface, a target called <code>multi-user.target</code>:

sudo ls -l /etc/systemd/system/default.target

/etc/systemd/system/default.target -> /usr/lib/systemd/system/multi-user.target



You can use a kernel boot parameter to override the default system target. See

systemd Units

systemd organizes the different types of resources it manages into units. The majority of units are configured in unit configuration files that enable you to configure these units according to your system needs. In addition to the files, you can also use systemd runtime commands to configure the units.

The following list describes some of the system units <code>systemd</code> enables you to manage on an Oracle Linux system:

Services

Service unit configuration files have the filename format service_name.service, for example sshd.service, crond.service, and httpd.service.

Service units start and control daemons and the processes of which the daemons consist.

The following example shows how you might start the systemd service unit for the Apache HTTP server, httpd.service:

sudo systemctl start httpd.service

Targets

Target unit configuration files have the filename format *target_name*.target, for example graphical.target.

Targets are similar to runlevels. A system reaches different targets during the boot process as resources get configured. For example, a system will reach <code>network-pre.target</code> before it reaches the target <code>network-online.target</code>.

Many target units have dependencies. For example, the activation of graphical.target (for a graphical session) will fail unless multi-user.target (for multi-user system) is also active.

File System Mount Points

Mount unit configuration files have the filename format $mount_point_name.mount$. Mount units enable you to mount filesystems at boot time. For example, you can run the following command to mount the temporary file system (tmpfs) on /tmp at boot time:

sudo systemctl enable tmp.mount

Devices

Device unit configuration files have the filename format <code>device_unit_name.device</code>. Device units are named after the <code>/sys</code> and <code>/dev</code> paths they control. For example, the <code>device /dev/sda5</code> is exposed in systemd as <code>dev-sda5.device</code>. Device units enable you to implement device-based activation.

Sockets

Socket unit configuration files have the filename format <code>socket_unit_name.socket</code>. Each "*.socket" file needs a corresponding "*.service" file to configure the service to start on incoming traffic on the socket.

Socket units enable you to implement socket-based activation.

See About System-State Targets.

About System-State Targets

By using system-state targets, you can control systemd so that it starts only the services that are required for a specific purpose. For example, you might wish to set the default target to multi-user.target on a production server so that it does not start the graphical user interface when it boots. In a case where you need to troubleshoot or perform diagnostics you might consider setting the target to rescue.target which only allows root to log onto the system and runs the minimum number of services.

Each run level defines the services that systemd stops or starts. As an example, systemd starts network services for multi-user.target and the X Window System for graphical.target, and stops both services for rescue.target.



Table 2-1 shows the commonly used system-state targets and the equivalent runlevel targets.

Table 2-1 System-State Targets and Equivalent Runlevel Targets

System-State Targets	Equivalent Runlevel Targets	Description
graphical.target	runlevel5.target	Set up a multi-user system with networking and display manager.
multi-user.target	runlevel2.target	Set up a non-graphical multi-
	runlevel3.target	user system with networking.
	runlevel4.target	
poweroff.target	runlevel0.target	Shut down and power off the system.
reboot.target	runlevel6.target	Shut down and reboot the system.
rescue.target	runlevel1.target	Set up a rescue shell.

Note that runlevel* targets are implemented as symbolic links.

For more information, see the systemd.target (5) manual page.

Displaying Default and Active System-State Targets

To display the default system-state target, use the systemctl get-default command:

```
sudo systemctl get-default
graphical.target
```

To display the currently active targets on a system, use the systemctl list-units -- type target command:

sudo systemctl list-units --type target [--all]

```
UNIT
basic.target loaded active active Basic System
cryptsetup.target loaded active active Local Encrypted Volumes
getty.target loaded active active Login Prompts
graphical.target loaded active active Graphical Interface
local-fs-pre.target loaded active active Local File Systems (Pre)
local-fs.target loaded active active Local File Systems
multi-user.target loaded active active Multi-User System
network-online.target loaded active active Network is Online
network-pre.target loaded active active Network (Pre)
network.target loaded active active Network
nfs-client.target loaded active active NFS client services
nss-user-lookup.target loaded active active User and Group Name Lookups
paths.target loaded active active Paths
remote-fs-pre.target loaded active active Remote File Systems (Pre)
remote-fs.target loaded active active Remote File Systems
rpc_pipefs.target loaded active active RPC Port Mapper
slices.target loaded active active Sockets
sound.target loaded active active Sockets
sound.target loaded active active Sound Card
```



```
sshd-keygen.target loaded active active sshd-keygen.target swap.target loaded active active Swap sysinit.target loaded active active System Initialization timers.target loaded active active Timers

LOAD = Reflects whether the unit definition was properly loaded.

ACTIVE = The high-level unit activation state, i.e. generalization of SUB. SUB = The low-level unit activation state, values depend on unit type.

24 loaded units listed. Pass --all to see loaded but inactive units, too. To show all installed unit files use 'systemctl list-unit-files'.
```

The output for a system with the graphical target active shows that this target depends on a number of other active targets, including network and sound to support networking and sound.

Use the --all option to include inactive targets in the displayed list.

For more information, see the systemctl(1) and systemd.target(5) manual pages.



Target is only one of multiple systemd types of units. To display all the types of units, use the following command:

```
sudo systemctl -t help
Available unit types:
service
mount
swap
socket
target
device
automount
timer
path
slice
scope
```

Changing Default and Active System-State Targets

Use the systemctl set-default command to change the default system-state target:

```
sudo systemctl set-default multi-user.target

Removed /etc/systemd/system/default.target.
Created symlink /etc/systemd/system/default.target → /usr/lib/systemd/system/multi-user.target
```





This command changes the target to which the default target is linked, but does not change the state of the system.

To change the currently active system target, use the <code>systemctl isolate</code> command, for example:

sudo systemctl isolate multi-user.target

For more information, see the systematl (1) manual page.

Shutting Down, Suspending, and Rebooting the System

Table 2-2 systemctl Commands for Shutting Down, Suspending, and Rebooting a System

systemctl Command	Description
systemctl halt	Halt the system.
systemctl hibernate	Put the system into hibernation.
systemctl hybrid-sleep	Put the system into hibernation and suspend its operation.
systemctl poweroff	Halt and power off the system.
systemctl reboot	Reboot the system.
systemctl suspend	Suspend the system.

For more information, see the systemctl(1) manual page.

Managing Services

Services in an Oracle Linux system are managed by the systemctl subcommand command.

systemctl supports a number of subcommands, for example enable, disable, stop, start, restart, reload, and status.

For more information, see the systematl(1) manual page.

Starting and Stopping Services

To start a service, use the systemctl start command:

sudo systemctl start sshd

To stop a service, use the systematl stop command:

sudo systemctl stop sshd





Changing the state of a service only lasts as long as the system remains at the same state. If you stop a service and then change the system-state target to one in which the service is configured to run (for example, by rebooting the system), the service restarts. Similarly, starting a service does not enable the service to start following a reboot. See Enabling and Disabling Services.

Enabling and Disabling Services

You can use the systemctl command to enable or disable a service from starting when the system boots, for example:

```
sudo systemctl enable httpd
```

Created symlink /etc/systemd/system/multi-user.target.wants/httpd.service \rightarrow /usr/lib/systemd/system/httpd.service.

The <code>enable</code> command enables a service by creating a symbolic link for the lowest-level system-state target at which the service should start. In the example, the command creates the symbolic link <code>httpd.service</code> for the <code>multi-user</code> target.

Disabling a service removes the symbolic link:

```
sudo systemctl disable httpd

Removed /etc/systemd/system/multi-user.target.wants/httpd.service.
```

To check whether a service is enabled, use is-enabled subcommand as shown in the following examples:

```
sudo systemctl is-enabled httpd
disabled
sudo systemctl is-enabled sshd
enabled
```

After running the systemctl disable command, the service can still be started or stopped by user accounts, scripts and other processes. However, if you need to ensure there is no possibility of the service being started inadvertently, for example, by a conflicting service, then you can use the systemctl mask command as follows:

```
sudo systemctl mask httpd
Created symlink from '/etc/systemd/system/multi-user.target.wants/httpd.service'
to '/dev/null'
```

The mask command sets the the service reference to /dev/null. If you try to start a service that has been masked, you will receive an error as shown in the following example:

```
sudo systemctl start httpd
```



```
Failed to start httpd.service: Unit is masked.
```

To re-link the service reference back to the matching service unit configuration file, use the systematl unmask command:

```
sudo systemctl unmask httpd
```

For more information, see the systemctl(1) manual page.

Displaying the Status of Services

To check whether a service is running, use the is-active subcommand. The output would either be *active*) or *inactive*, as shown in the following examples:

```
sudo systemctl is-active httpd
active
systemctl is-active sshd
inactive
```

The status subcommand provides a detailed summary of the status of a service, including a tree that displays all of the tasks in the *control group* (CGroup) that the service implements:

```
sudo systemctl status httpd
httpd.service - The Apache HTTP Server
   Loaded: loaded (/usr/lib/systemd/system/httpd.service; enabled; vendor preset:
  Active: active (running) since ...
    Docs: man:httpd.service(8)
Main PID: 11832 (httpd)
  Status: "Started, listening on: port 80"
   Tasks: 213 (limit: 26213)
   Memory: 32.5M
   CGroup: /system.slice/httpd.service
           -11832 /usr/sbin/httpd -DFOREGROUND
           -11833 /usr/sbin/httpd -DFOREGROUND
           -11834 /usr/sbin/httpd -DFOREGROUND
           -11835 /usr/sbin/httpd -DFOREGROUND
           L11836 /usr/sbin/httpd -DFOREGROUND
Jul 17 00:14:32 Unknown systemd[1]: Starting The Apache HTTP Server...
Jul 17 00:14:32 Unknown httpd[11832]: Server configured, listening on: port 80
Jul 17 00:14:32 Unknown systemd[1]: Started The Apache HTTP Server.
```

A *cgroup* is a collection of processes that are bound together so that you can control their access to system resources. In the example, the cgroup for the httpd service is httpd.service, which is in the system *slice*.

Slices divide the cgroups on a system into different categories. To display the slice and cgroup hierarchy, use the systemd-cgls command:

```
sudo systemd-cgls
Control group /:
-.slice
|-user.slice
| Luser-1000.slice
```



```
-user@1000.service
    └init.scope
      -6488 /usr/lib/systemd/systemd --user
      └-6492 (sd-pam)
    -session-7.scope
    -6484 sshd: root [priv]
    -6498 sshd: root@pts/0
    -6499 -bash
     -6524 sudo systemd-cgls
     -6526 systemd-cgls
    └6527 less
-init.scope
-1 /usr/lib/systemd/systemd --switched-root --system --deserialize 16
-system.slice
 -rngd.service
 └1266 /sbin/rngd -f --fill-watermark=0
 -irqbalance.service
 └1247 /usr/sbin/irgbalance --foreground
 -libstoragemgmt.service
 └1201 /usr/bin/lsmd -d
 -systemd-udevd.service
 └1060 /usr/lib/systemd/systemd-udevd
  -polkit.service
  └1241 /usr/lib/polkit-1/polkitd --no-debug
  -chronyd.service
  └1249 /usr/sbin/chronyd
 -auditd.service
  ⊢1152 /sbin/auditd
  L1154 /usr/sbin/sedispatch
 -tuned.service
 └-1382 /usr/libexec/platform-python -Es /usr/sbin/tuned -l -P
 -systemd-journald.service
  └1027 /usr/lib/systemd/systemd-journald
  -atd.service
  └1812 /usr/sbin/atd -f
  -sshd.service
 └1781 /usr/sbin/sshd
```

The system.slice contains services and other system processes. user.slice contains user processes, which run within transient cgroups called *scopes*. In the example, the processes for the user with ID 1000 are running in the scope session-7.scope under the slice /user.slice/user-1000.slice.

You can use the systemctl command to limit the CPU, I/O, memory, and other resources that are available to the processes in service and scope cgroups. See Controlling Access to System Resources.

For more information, see the systemctl(1) and systemd-cqls(1) manual pages.

Controlling Access to System Resources

Use the ${\tt systemctl}$ command to control a cgroup's access to system resources, for example:

```
sudo systemctl [--runtime] set-property httpd CPUShares=512 MemoryLimit=1G
```

CPUShare controls access to CPU resources. As the default value is 1024, a value of 512 halves the access that the processes in the cgroup have to CPU time. Similarly, MemoryLimit controls the maximum amount of memory that the cgroup can use.

Note:

You do not need to specify the .service extension to the name of a service.

If you specify the --runtime option, the setting does not persist across system reboots.

Alternatively, you can change the resource settings for a service under the <code>[Service]</code> heading in the service's configuration file in <code>/usr/lib/systemd/system</code>. After editing the file, make <code>systemd</code> reload its configuration files and then restart the service:

```
sudo systemctl daemon-reload
sudo systemctl restart service
```

You can run general commands within scopes and use systemctl to control the access that these transient cgroups have to system resources. To run a command within in a scope, use the systemd-run command:

```
sudo systemd-run --scope --unit=group name [--slice=slice name]
```

If you do not want to create the group under the default system slice, you can specify another slice or the name of a new slice. The following example runs a command named mymonitor in mymon.scope under myslice.slice:

```
sudo systemd-run --scope --unit=mymon --slice=myslice mymonitor
Running as unit mymon.scope.
```



If you do not specify the --scope option, the control group is a created as a service rather than as a scope.

You can then use systemctl to control the access that a scope has to system resources in the same way as for a service. However, unlike a service, you must specify the .scope extension, for example:

```
sudo systemctl --runtime set-property mymon.scope CPUShares=256
```

For more information see the systemctl(1), systemd-cgls(1), and systemd.resource-control(5) manual pages.

Running systemctl on a Remote System

If the sshd service is running on a remote Oracle Linux system, specify the -H option with the systemctl command to control the system remotely, for example:

```
sudo systemctl -H root@10.0.0.2 status sshd
root@10.0.0.2's password: password
sshd.service - OpenSSH server daemon
  Loaded: loaded (/usr/lib/systemd/system/sshd.service; enabled)
```



```
Active: active (running) since ...

Process: 1498 ExecStartPre=/usr/sbin/sshd-keygen (code=exited, status=0/SUCCESS)

Main PID: 1524 (sshd)

CGroup: /system.slice/sshd.service
```

For more information see the systemctl(1) manual page.

Modifying systemd Service Unit Files

To change the configuration of systemd services, copy all of the files with .service, .target, .mount and .socket extensions from /usr/lib/systemd/system to /etc/systemd/system.

Once you have copied the files you can edit the versions in <code>/etc/systemd/system</code>. The copies of the files in <code>/etc/systemd/system</code> take precedence over the version in <code>/usr/lib/systemd/system</code>, and are not overwritten when you update a package that touches files in <code>/usr/lib/systemd/system</code>.

To revert to the default systemd configuration for a particular service, either rename or delete the modified copies in /etc/systemd/system.

The following sections look at the different parts of a service unit file that you might wish to edit and customize for your own system.

About Service Unit Files

Services run based on their corresponding service unit files. A service unit file typically contains the following sections, with each section having its respective defined options that determine how a specific service runs:

[Unit]

The [Unit] section contains information about the service.

[UnitType]:

The [UnitType] section contains options that are specific to the unit type of the file. For example, in a service unit file this section is titled [Service] and contains options that are specific to units of the service type, such as ExecStart or StandardOutput. Only those unit types that offer options specific to their type have such a section.

[Install]

The [Install] section contains installation information for the specific unit. The information in this section is used by the systematl enable and systematl disable commands.

A service unit file might contain the following configurations for a service.

```
[Unit]
Description=A test service used to develop a service unit file template
[Service]
Type=simple
StandardOutput=journal
ExecStart=/usr/lib/systemd/helloworld.sh
```



```
[Install]
WantedBy=default.target
```

The following section, Configurable Options in Service Unit Files, describes some of the commonly used configured options available under each section. A complete list is also available in the systemd.service(5) and systemd.unit(5) manual pages.

Configurable Options in Service Unit Files

Each of the following lists deals with a separate section of the service unit file.

Description of Options Under [Unit] Section

The following list provides a general overview of the commonly used configurable options available in the [Unit] section of service unit file:

Description

The Description option provides information about the service. The information is displayed when you run the systemctl status command on the unit.

Documentation

The Documentation option contains a space-separated list of URIs referencing documentation for this unit or its configuration.

After

The After option configures the unit to only run after the units listed in the option finish starting up.

In the following example, if the file *var3*.service has the following entry then it is only started after units *var1*.service and *var2*.service have started:

After=var1.service var2.service

Requires

The Requires option configures a unit to have requirement dependencies on other units. If a unit is activated, those listed in its Requires option are also activated.

Wants

The Wants option Is a less stringent version of the Requires option. For example, a given unit can be activated even if one of those listed in its Wants option fails to start.

Description of Options Under [Service] Section

This following list gives a general overview of the commonly used configurable options available in the [Service] section of a service unit file.

Type

The Type option configures the process start-up type for the service unit.

By default, this parameter's value is simple, which indicates that the service's main process is that which is started by the ExecStart parameter.

Typically, if a service's type is simple, then the definition can be omitted from the file.

StandardOutput

The StandardOutput option enables you configure the how the service's events are logged. For example, consider a service unit file has the following entry:

StandardOutput=journal



In the example, the value journal indicates that the events are recorded in the journal, which can be viewed by using the journalctl command.

ExecStart

The ExecStart option specifies the full path and command that starts the service, for example, /usr/bin/npm start.

ExecStop

The ExecStop option specifies the commands to run to stop the service started via ExecStart.

ExecReload

The ExecReload option specifies the commands to run to trigger a configuration reload in the service.

Restart

The Restart option configures whether the service is to be restarted when the service process exits, is terminated, or a timeout is reached.



This option does not apply when the process is stopped cleanly by a systemd operation, for example a systemctl stop or systemctl restart. In these cases, the service is not restarted by this configuration option.

RemainAfterExit

The RemainAfterExit option is a Boolean value that configures whether the service is to be considered active even when all of its processes have exited. The default value is no.

Description of Options Under [Install] Section

This following list gives a general overview of the commonly used configurable options available in the [Install] section of service unit file.

Alias

The Alias option is a space-separated list of additional names for a unit.

At installation time, systemctl enable will create symlinks from these names to the unit filename.

Aliases are only effective when the unit is enabled.

RequiredBy

The RequiredBy option enables you to configure the service to be required by other units

For example, consider a unit file var1.service that has the following configuration added to it:

RequiredBy=var2.service var3.service

When var1.service is enabled, both var2.service and var3.service are given a Requires dependency upon var1.service. This dependency is defined by a symbolic link that is created in the ".requires" folder of each dependent service (var2.service and var3.service) that points to the var1.service system unit file.



WantedBy

The WantedBy option enables you to specify a list of units that are to be given a wants dependency upon the service whose file you are editing.

For example, consider a unit file *var1*.service that has the following configuration added to it:

WantedBy=var2.service var3.service

When var1.service is enabled, both var2.service and var3.service are given a Wants dependency upon var1.service. This dependency is defined by a symbolic link that is created in the ".wants" folder of each dependent service (var2.service and var3.service) that points to the system unit file for var1.service.

Also

The Also option lists additional units to install or remove when the unit is installed or removed.

DefaultInstance

The DefaultInstance option applies to template unit files only.

Template unit files allow creation of multiple units from a single configuration file. The <code>DefaultInstance</code> option specifies the instance for which the unit is enabled if the template is enabled without any explicitly set instance.

Creating a User-Based systemd Service

In addition to the system-wide systemd files, systemd enables you to create user-based services that you can run from a user level without requiring root access and privileges. These user-based services are under your control and are configurable independent of system services.

The following are some distinguishing features of user-based systemd services:

- User-based systemd services are linked with a specific user account.
- They are generally created under the associated user's home directory in \$HOME/.config/systemd/user/.
- Once enabled, the services will start when the associated user logs on (unlike systemwide systemd services which, once enabled, start when the system boots).

The feature is particularly useful when creating podman container services. For more information about podman, see Oracle Linux: Podman User's Guide.

To create a user based service:

- 1. Create your service's unit file in the ~/.config/systemd/user directory, for example:
 - touch ~/.config/systemd/user/myservice.service
- 2. Open the unit file and specify the values to the options you want to use, such as Description, ExecStart, WantedBy, and so on.
 - For reference, see Configurable Options in Service Unit Files and the systemd.service(5) and systemd.unit(5) manual pages.
- 3. Enable the service to start automatically when you log in.

```
sudo systemctl --user enable myservice.service
```



Note:

When you log out, the service is generally stopped unless the root user has enabled processes to continue to run for your user. See Use systemd on Oracle Linux for more information.

4. Start the service.

sudo systemctl --user start myservice.service

5. Verify that the service is running.

sudo systemctl --user status myservice.service



Configuring System Settings

This chapter describes the files and virtual file systems that you can use to change the configuration settings for your system.

Also see Configure System Settings on Oracle Linux for a hands-on tutorial on how to configure system settings as described in this chapter.

About the /etc/sysconfig Files

The /etc/sysconfig directory contains files that control your system's configuration. The contents of this directory depend on the packages that you have installed on your system.

Some of the files that you might find in the /etc/sysconfig directory include the following:

atd

Specifies additional command line arguments for the atd daemon.

autofs

Defines custom options for automatically mounting devices and controlling the operation of the automounter.

crond

Passes arguments to the crond daemon at boot time.

chronyd

Passes arguments to the chronyd daemon used for NTP services at boot time.

firewalld

Passes arguments to the firewall daemon (firewalld) at boot time.

grub

Specifies default settings for the GRUB 2 bootloader. This file is a symbolic link to /etc/default/grub. For more information, see About the GRUB 2 Bootloader.

named

Passes arguments to the name service daemon at boot time. The named daemon is a Domain Name System (DNS) server that is part of the Berkeley Internet Name Domain (BIND) distribution. This server maintains a table that associates host names with IP addresses on the network.

samba

Passes arguments to the smbd, nmbd, and winbindd daemons at boot time to support file-sharing connectivity for Windows clients, NetBIOS-over-IP naming service, and connection management to domain controllers.

selinux

Controls the state of SELinux on the system. This file is a symbolic link to /etc/selinux/config.

For more information, see Oracle Linux: Administering SELinux.

snapper

Defines a list of btrfs file systems and thinly-provisioned LVM volumes whose contents can be recorded as snapshots by the snapper utility.

For more information, see Oracle Linux 9: Managing Local File Systems.

sysstat

Configures logging parameters for system activity data collector utilities such as sar.

For more information, see /usr/share/doc/initscripts*/sysconfig.txt.

About the /proc Virtual File System

The files in the /proc directory hierarchy contain information about your system hardware and the processes that are running on the system. You can change the configuration of the kernel by writing to certain files that have write permission.

Files that are under the /proc directory are virtual files that the kernel creates on demand to present a browsable view of the underlying data structures and system information. As such, /proc is an example of a virtual file system. Most virtual files are listed as zero bytes in size, but they contain large amount of information when viewed.

Virtual files such as /proc/interrupts, /proc/meminfo, /proc/mounts, and /proc/partitions provide a view of the system's hardware. Other files, such as /proc/filesystems and the files under /proc/sys, provide information about the system's configuration and allow this configuration to be modified.

Files that contain information about related topics are grouped into virtual directories. A separate directory exists in the /proc directory for each process that is currently running on the system; the directory's name corresponds to the numeric process ID. For example, /proc/1 corresponds to the systemd process that has a PID of 1.

To examine virtual files, you can use commands such as cat, less, and view, as shown in the following example:

cat /proc/cpuinfo

```
processor : 0
              : GenuineIntel
cpu family model
vendor id
              : 42
model name
              : Intel(R) Core(TM) i5-2520M CPU @ 2.50GHz
stepping
               : 7
              : 2393.714
cpu MHz
cache size physical id
              : 6144 KB
              : 0
siblings
               : 2
core id
              : 0
cpu cores
              : 2
apicid
initial apicid : 0
              : yes
fpu exception : yes
cpuid level
              : 5
wp
               : yes
```



For files that contain nonhuman-readable content, you can use utilities such as lspci, free, top, and sysctl to access information. For example, the lspci command lists PCI devices on a system:

sudo lspci

Virtual Files and Directories Under /proc

The following table describes the most useful virtual files and directories under the /proc directory hierarchy.



Table 3-1 Useful Virtual Files and Directories Under the /proc Directory

Virtual File or Directory	Description
PID (Directory)	Provides information about the process with the process ID (<i>PID</i>). The directory's owner and group is same as the process's. Useful files under the directory include:
	cmdline Command path.
	cwd Symbolic link to the process's current working directory.
	environ Environment variables.
	exe Symbolic link to the command executable.
	fd/N File descriptors.
	maps Memory maps to executable and library files.
	root Symbolic link to the effective root directory for the process.
	stack The contents of the kernel stack.
	status Run state and memory usage.
buddyinfo	Provides information for diagnosing memory fragmentation.
bus (directory)	Contains information about the various buses (such as pci and usb) that are available on the system. You can use commands such as lspci, lspcmcia, an lsusb to display information for such devices.
cgroups	Provides information about the resource control groups that are in use on the system.
cmdline	Lists parameters passed to the kernel at boot time.
cpuinfo	Provides information about the system's CPUs.



Table 3-1 (Cont.) Useful Virtual Files and Directories Under the /proc Directory

Virtual File or Directory	Description
crypto	Provides information about all installed cryptographic cyphers.
devices	Lists the names and major device numbers of all currently configured characters and block devices.
dma	Lists the direct memory access (DMA) channels that are currently in use.
driver (directory)	Contains information about drivers used by the kernel, such as those for non-volatile RAM (nvram), the real-time clock (rtc), and memory allocation for sound (snd-page-alloc).
execdomains	Lists the execution domains for binaries that the Oracle Linux kernel supports.
filesystems	Lists the file system types that the kernel supports. Entries marked with nodev are not in use.
fs (directory)	Contains information about mounted file systems, organized by file system type.
interrupts	Records the number of interrupts per interrupt request queue (IRQ) for each CPU since system startup.
iomem	Lists the system memory map for each physical device.
ioports	Lists the range of I/O port addresses that the kernel uses with devices.
irq(directory)	Contains information about each IRQ. You can configure the affinity between each IRQ and the system CPUs.
kcore	Presents the system's physical memory in core file format that you can examine using a debugger such as crash or gdb. This file is not human-readable.
kmsg	Records kernel-generated messages, which are picked up by programs such as dmesg.
loadavg	Displays the system load averages (number of queued processes) for the past 1, 5, and 15 minutes, the number of running processes, the total number of processes, and the PID of the process that is running.



Table 3-1 (Cont.) Useful Virtual Files and Directories Under the /proc Directory

Virtual File or Directory	Description
locks	Displays information about the file locks that the kernel is currently holding on behalf of processes. The information provided includes:
	 lock class (FLOCK or POSIX)
	 lock type (ADVISORY or MANDATORY)
	 access type (READ or WRITE)
	process IDmajor device, minor device, and inode numbers
	 bounds of the locked region
mdstat	Lists information about multiple-disk RAID devices.
meminfo	Reports the system's usage of memory in more detail than is available using the free or top commands.
modules	Displays information about the modules that are currently loaded into the kernel. The lsmod command formats and displays the same information, excluding the kernel memory offset of a module.
mounts	Lists information about all mounted file systems.
net (directory)	Provides information about networking protocol, parameters, and statistics. Each directory and virtual file describes aspects of the configuration of the system's network.
partitions	Lists the major and minor device numbers, number of blocks, and name of partitions mounted by the system.
scsi/device_info	Provides information about supported SCSI devices.
scsi/scsi and	Provide information about configured SCSI
scsi/sg/*	devices, including vendor, model, channel, ID, and LUN data .
self	Symbolic link to the process that is examining /proc.
slabinfo	Provides detailed information about slab memory usage.
softirgs	Displays information about software interrupts (<i>softirqs</i>). A softirq is similar to a hardware interrupt (<i>hardirq</i>) and allow the kernel to perform asynchronous processing that would take too long during a hardware interrupt.



Table 3-1 (Cont.) Useful Virtual Files and Directories Under the /proc Directory

Virtual File or Directory	Description
stat	Records information about the system since it was started, including:
	cpu Total CPU time (measured in <i>jiffies</i>) spent in user mode, low-priority user mode, system mode, idle, waiting for I/O, handling hardirq events, and handling softirq events.
	cpuN Times for CPU <i>N</i> .
swaps	Provides information about swap devices. The units of size and usage are kilobytes.
sys (directory)	Provides information about the system and also allows you to enable, disable, or modify kernel features. You can write new settings to any file that has write permission. See Modifying Kernel Parameters.
	The following subdirectory hierarchies of proc/sys contain virtual files, some of whose values you can usefully alter:
	dev Device parameters.
	fs File system parameters.
	kernel Kernel configuration parameters.
	net Networking parameters.
sysvipc (directory)	Provides information about the usage of System V Interprocess Communication (IPC) resources for messages (msg), semaphores (sem), and shared memory (shm).
tty (directory)	Provides information about the available and currently used terminal devices on the system. The drivers virtual file lists the devices that are currently configured.
vmstat	Provides information about virtual memory usage.

For more information, see the proc(5) manual page.



Modifying Kernel Parameters

Some virtual files under /proc, and especially under /proc/sys, are writable. You can adjust settings in the kernel through these files. For example, to change the hostname, you would revise the /proc/sys/kernel/hostname file as follows:

```
echo www.mydomain.com > /proc/sys/kernel/hostname
```

Other files take binary or boolean values, such as the setting of IP forwarding, which is defined in /proc/sys/net/ipv4/ip forward:

```
cat /proc/sys/net/ipv4/ip_forward
0
echo 1 > /proc/sys/net/ipv4/ip_forward
cat /proc/sys/net/ipv4/ip_forward
1
```

You can use the sysctl command to view or modify values under the /proc/sys directory.



Even root cannot bypass the file access permissions of virtual file entries under /proc. If you attempt to change the value of a read-only entry such as /proc/partitions, there is no kernel code to service the write() system call.

To display all of the current kernel settings, use the following command:

```
sysctl -a
```

```
kernel.sched_child_runs_first = 0
kernel.sched_min_granularity_ns = 2000000
kernel.sched_latency_ns = 10000000
kernel.sched_wakeup_granularity_ns = 2000000
kernel.sched_shares_ratelimit = 500000
```

Note:

The delimiter character in the name of a setting is a period (.) rather than a slash (/) in a path relative to /proc/sys, such as net.ipv4.ip_forward. This setting represents net/ipv4/ip_forward. As another example, kernel.msgmax represents kernel/msgmax.

, ,

To display an individual setting, specify its name as the argument to sysctl:

```
sysctl net.ipv4.ip_forward
net.ipv4.ip forward = 0
```

To change the value of a setting, use the following command format:

```
sysctl -w net.ipv4.ip_forward=1
net.ipv4.ip forward = 1
```

Changes that you make in this way remain in force only until the system is rebooted. To make configuration changes persist after the system is rebooted, you must add them to the /etc/sysctl.d directory as a configuration file. Any changes that you make to the files in this directory take effect when the system reboots or if you run the sysctl --system command, for example:

```
echo 'net.ipv4.ip_forward=1' > /etc/sysctl.d/ip_forward.confgrep -r ip_forward /etc/
sysctl.d
/etc/sysctl.d/ip forward.conf:net.ipv4.ip forward=1
sysctl net.ipv4.ip forward
net.ipv4.ip forward = 0
sysctl --system
* Applying /usr/lib/sysctl.d/00-system.conf ...
net.bridge.bridge-nf-call-ip6tables = 0
net.bridge.bridge-nf-call-iptables = 0
net.bridge.bridge-nf-call-arptables = 0
* Applying /usr/lib/sysctl.d/50-default.conf ...
kernel.sysrq = 16
kernel.core uses pid = 1
net.ipv4.conf.default.rp filter = 1
net.ipv4.conf.all.rp filter = 1
net.ipv4.conf.default.accept source route = 0
net.ipv4.conf.all.accept source route = 0
net.ipv4.conf.default.promote secondaries = 1
net.ipv4.conf.all.promote secondaries = 1
fs.protected hardlinks = \overline{1}
fs.protected symlinks = 1
* Applying /etc/sysctl.d/99-sysctl.conf ...
* Applying /etc/sysctl.d/ip_forward.conf ...
net.ipv4.ip forward = 1
* Applying /etc/sysctl.conf ...
sysctl net.ipv4.ip forward
net.ipv4.ip forward = 1
```

For more information, see the sysctl(8) and sysctl.d(5) manual pages.

Parameters That Control System Performance

The following parameters control various aspects of system performance:

fs.file-max

Specifies the maximum number of open files for all processes. Increase the value of this parameter if you see messages about running out of file handles.

net.core.netdev max backlog

Specifies the size of the receiver backlog queue, which is used if an interface receives packets faster than the kernel can process them. If this queue is too small, packets are lost at the receiver, rather than on the network.

net.core.rmem max

Specifies the maximum read socket buffer size. To minimize network packet loss, this buffer must be large enough to handle incoming network packets.



net.core.wmem max

Specifies the maximum write socket buffer size. To minimize network packet loss, this buffer must be large enough to handle outgoing network packets.

net.ipv4.tcp available congestion control

Displays the TCP congestion avoidance algorithms that are available for use. Use the modprobe command if you need to load additional modules such as tcp_htcp to implement the htcp algorithm.

net.ipv4.tcp congestion control

Specifies which TCP congestion avoidance algorithm is used.

net.ipv4.tcp max syn backlog

Specifies the number of outstanding SYN requests that are allowed. Increase the value of this parameter if you see <code>synflood</code> warnings in your logs, and investigation shows that they are occurring because the server is overloaded by legitimate connection attempts.

net.ipv4.tcp rmem

Specifies minimum, default, and maximum receive buffer sizes that are used for a TCP socket. The maximum value cannot be larger than net.core.rmem max.

net.ipv4.tcp wmem

Specifies minimum, default, and maximum send buffer sizes that are used for a TCP socket. The maximum value cannot be larger than net.core.wmem max.

vm.swappiness

Specifies how likely the kernel is to write loaded pages to swap rather than drop pages from the system page cache. When set to 0, swapping only occurs to avoid an out of memory condition. When set to 100, the kernel swaps aggressively. For a desktop system, setting a lower value can improve system responsiveness by decreasing latency. The default value is 60.



Caution:

This parameter is intended for use with laptop computers to reduce power consumption by the hard disk. Do not adjust this value on server systems.

Parameters That Control Kernel Panics

The following parameters control the circumstances under which a kernel panic can occur:

kernel.hung task panic

If set to 1, the kernel panics if any kernel or user thread sleeps in the TASK_UNINTERRUPTIBLE state (*D state*) for more than

kernel.hung_task_timeout_secs seconds. A process remains in D state while waiting for I/O to complete. You cannot kill or interrupt a process in this state. The default value is 0, which disables the panic.





Tip:

To diagnose a hung thread, you can examine /proc/PID/stack, which displays the kernel stack for both kernel and user threads.

kernel.hung_task_timeout_secs

Specifies how long a user or kernel thread can remain in D state before a warning message is generated or the kernel panics, if the value of kernel.hung_task_panic is 1. The default value is 120 seconds. A value of 0 disables the timeout.

kernel.nmi watchdog

If set to 1 (default), enables the non-maskable interrupt (NMI) watchdog thread in the kernel. If you want to use the NMI switch or the OProfile system profiler to generate an undefined NMI, set the value of kernel.nmi watchdog to 0.

kernel.panic

Specifies the number of seconds after a panic before a system will automatically reset itself. If the value is 0, the system hangs, which allows you to collect detailed information about the panic for troubleshooting. This is the default value.

To enable automatic reset, set a non-zero value. If you require a memory image (vmcore), allow enough time for Kdump to create this image. The suggested value is 30 seconds, although large systems will require a longer time.

kernel.panic on io nmi

If set to 0 (default), the system tries to continue operations if the kernel detects an I/O channel check (IOCHK) NMI that usually indicates a uncorrectable hardware error. If set to 1, the system panics.

kernel.panic on oops

If set to 0, the system tries to continue operations if the kernel encounters an oops or BUG condition. If set to 1 (default), the system delays a few seconds to give the kernel log daemon, klogd, time to record the oops output before the panic occurs.

In an OCFS2 cluster, set the value to 1 to specify that a system must panic if a kernel oops occurs. If a kernel thread required for cluster operation crashes, the system must reset itself. Otherwise, another node might not be able to tell whether a node is slow to respond or unable to respond, causing cluster operations to hang.

kernel.panic on unrecovered nmi

If set to 0 (default), the system tries to continue operations if the kernel detects an NMI that usually indicates an uncorrectable parity or ECC memory error. If set to 1, the system panics.

kernel.softlockup panic

If set to 0 (default), the system tries to continue operations if the kernel detects a *soft-lockup* error that causes the NMI watchdog thread to fail to update its time stamp for more than twice the value of kernel.watchdog thresh seconds. If set to 1, the system panics.

kernel.unknown nmi panic

If set to 1, the system panics if the kernel detects an undefined NMI. You would usually generate an undefined NMI by manually pressing an NMI switch. As the NMI watchdog thread also uses the undefined NMI, set the value of $kernel.unknown_nmi_panic$ to 0 if you set kernel.nmi watchdog to 1.



kernel.watchdog thresh

Specifies the interval between generating an NMI performance monitoring interrupt that the kernel uses to check for *hard-lockup* and *soft-lockup* errors. A hard-lockup error is assumed if a CPU is unresponsive to the interrupt for more than kernel.watchdog_thresh seconds. The default value is 10 seconds. A value of 0 disables the detection of lockup errors.

vm.panic_on_oom

If set to 0 (default), the kernel's OOM-killer scans through the entire task list and attempts to kill a memory-hogging process to avoid a panic. If set to 1, the kernel panics but can survive under certain conditions. If a process limits allocations to certain nodes by using memory policies or cpusets, and those nodes reach memory exhaustion status, the OOM-killer can kill one process. No panic occurs in this case because other nodes' memory might be free and the system as a whole might not yet be out of memory. If set to 2, the kernel always panics when an OOM condition occurs. Settings of 1 and 2 are for intended for use with clusters, depending on your preferred failover policy.

About the /sys Virtual File System

In addition to the /proc file system, the kernel exports information to the /sys virtual file system (sysfs). Programs such as the dynamic device manager (udev), use /sys to access device and device driver information.



/sys exposes kernel data structures and control points, which implies that it might contain circular references, where a directory links to an ancestor directory. As a result, a find command used on /sys might never terminate.

Virtual Directories Under the /sys Directory

The following table describes some useful virtual directories under the $/ \mathtt{sys}$ directory hierarchy.

Table 3-2 Virtual Directories Under /sys

Virtual Directory	Description
block	Contains subdirectories for block devices. For example: /sys/block/sda.
bus	Contains subdirectories for each supported physical bus type, such as pci, pcmcia, scsi, or usb. Under each bus type, the devices directory lists discovered devices, and the drivers directory contains directories for each device driver.
class	Contains subdirectories for every class of device that is registered with the kernel.



Table 3-2 (Cont.) Virtual Directories Under /sys

Virtual Directory	Description
dev	Contains the char/ and block/ directories. Inside these two directories there are symlinks named <major>:<minor>. These symlinks point to the sysfs directory for the given device. The /sys/dev directory provides a quick way to look up the sysfs interface for a device from the result of the stat (2) operation.</minor></major>
devices	Contains the global device hierarchy of all devices on the system. The platform directory contains peripheral devices such as device controllers that are specific to a particular platform. The system directory contains non-peripheral devices such as CPUs and APICs. The virtual directory contains virtual and pseudo devices. See Managing System Devices.
firmware	Contains subdirectories for firmware objects.
fs	Contains subdirectories for file system objects.
kernel	Contains subdirectories for additional kernel objects
module	Contains subdirectories for each module loaded into the kernel. You can alter some parameter values for loaded modules. See About Module Parameters.
power	Contains attributes that control the system's power state.

For more information, see https://www.kernel.org/doc/Documentation/filesystems/sysfs.txt.

Configuring System Date and Time Settings

System time is based on the POSIX time standard, where time is measured as the number of seconds that have elapsed since 00:00:00 Coordinated Universal Time (UTC), Thursday, 1 January 1970. A day is defined as 86400 seconds and leap seconds are subtracted automatically.

Date and time representation on a system can be set to match a specific timezone. To list all of the available timezones, run:

timedatectl list-timezones

To set the system timezone to match a value returned from the available timezones, you can run:

timedatectl set-timezone America/Los_Angeles

Substitute America/Los_Angeles with a valid timezone entry.



This command sets a symbolic link from /etc/localtime to point to the appropriate zone information file in /usr/share/zoneinfo/. The setting takes effect immediately. Some long running processes that might use /etc/localtime to detect the current system timezone, may not detect a subsequent change in system timezone until the process is restarted.

Note that timezones are largely used for display purposes or to handle user input. Changing timezone does not change the time for the system clock. You can change the presentation for system time in any console by setting the ${\tt TZ}$ environment variable. For example, to see the current time in Tokyo, you can run:

```
TZ="Asia/Tokyo" date
```

You can check your system's current date and time configuration by running the timedatectl command on its own:

timedatectl

```
Local time: Wed 2021-07-17 00:50:58
EDT

Universal time: Wed 2021-07-17 04:50:58
UTC

RTC time: Wed 2021-07-17
04:50:55

Time zone: America/New_York (EDT, -0400)

System clock synchronized: yes

NTP service: active

RTC in local TZ: no
```

To set system time manually, you can use the timedatectl set-time command:

```
timedatectl set-time "2021-07-17 01:59:59"
```

This command sets the current system time based on the time specified assuming the currently set system timezone. The command also updates the system Real Time Clock (RTC).



Tip:

See Learn How to Localize Your Installation on Oracle Linux for a hands-on tutorial that describes how to use tools to configure system parameters like date, time and locale.



Consider configuring your system to use network time synchronization for accurate time-keeping. This can be particularly important when setting up high-availability or when using network-based file systems.



Tip:

SeeConfigure Chrony on Oracle Linux for a hands-on tutorial on setting up and configuring the chronyd service.

If you configure an NTP service, enable NTP by running the following command:

timedatectl set-ntp true

This command enables and starts the chronyd service, if available.



4

Managing System Devices

This chapter describes how the system uses device files and how the Udev device manager dynamically creates or removes device node files.

About Device Files

The /dev directory contains device files or device nodes that provide access to peripheral devices such as hard disks, to resources on peripheral devices such as disk partitions, and pseudo devices such as a random number generator.

The /dev directory has several subdirectory hierarchies, each of which holds device files that relate to a certain type of device. However, the contents of these subdirectories are implemented as symbolic links to corresponding files in /dev. Thus, the files can be accessed either through the linked file in /dev or the corresponding file in the subdirectory.

Using the ls-1 /dev command lists files, some of which are flagged as being either type b (for *block*) or type c (for *character*). These devices have an associated pair of numbers that identify the device to the system.

```
ls -1 /dev
 total 0
 crw-r--r-. 1 root root 10, 235 Aug 20 08:36 autofs

      drwxr-xr-x.
      2 root root
      240 Sep 20 07:37 block

      drwxr-xr-x.
      2 root root
      100 Aug 20 08:36 bsg

      drwxr-xr-x.
      3 root root
      60 Nov 4 2019 bus

      lrwxrwxrwx.
      1 root root
      3 Aug 20 08:36 cdrom ->

      drwxr-xr-x.
      2 root root
      2720 Sep 20 07:37 char

      crw-----.
      1 root root
      5, 1 Aug 20 08:36 console

      lrwxrwxrwx.
      1 root root
      11 Aug 20 08:36 core ->

                                                                     3 Aug 20 08:36 cdrom -> sr0
                                                        11 Aug 20 08:36 core -> /proc/kcore
 lrwxrwxrwx. 1 root root
                                                                  60 Nov 4 2019 cpu
drwxr-xr-x. 3 root root
 crw-----. 1 root root 10, 62 Aug 20 08:36 cpu dma latency
                                                       140 Aug 20 08:36 disk
 drwxr-xr-x. 7 root root
brw-rw---. 1 root disk 253, 0 Aug 20 08:36 dm-0
brw-rw---. 1 root disk 253, 1 Aug 20 08:36 dm-1
brw-rw---. 1 root disk 253, 2 Aug 20 08:36 dm-2
lrwxrwxrwx. 1 root root
                                                                  13 Aug 20 08:36 fd -> /proc/self/fd

      crw-rw-rw-.
      1 root root
      10, 229 Aug 20 08:36 fuse

      crw------
      1 root root
      10, 228 Aug 20 08:36 hpet

      drwxr-xr-x.
      2 root root
      0 Aug 20 08:36 hugepages

      crw------
      1 root root
      12 Aug 20 08:36 initctl -> /r

      drwxr-xr-x.
      3 root root
      220 Aug 20 08:36 input

      crw-r-----
      1 root root
      1, 11 Aug 20 08:36 kmsg

      lrwxrwxrwx.
      1 root root
      28 Aug 20 08:36 log -> /run/s

      brw-rw----
      1 root disk
      7, 0 Sep 23 01:28 loop0

      crw-rw-----
      1 root disk
      10, 237 Sep 20 07:37 loop-control

      drwxr-xr-x.
      2 root root
      120 Aug 20 08:36 mapper

      crw------
      1 root root
      120 Aug 20 08:36 mcelog

                                                                  12 Aug 20 08:36 initctl -> /run/initctl
                                                                     28 Aug 20 08:36 log -> /run/systemd/journal/dev-log
 crw----. 1 root root
                                                         10, 227 Aug 20 08:36 mcelog
 crw-r---. 1 root kmem
                                                        1, 1 Aug 20 08:36 mem
 crw-----. 1 root root 10, 59 Aug 20 08:36 memory bandwidth
                                                                   40 Nov 4 2019 mqueue
 drwxrwxrwt. 2 root root
```

```
      drwxr-xr-x.
      2 root root
      60 Aug 20 08:36 net

      crw------.
      1 root root
      10, 61 Aug 20 08:36 network_latency

      crw------.
      1 root root
      10, 60 Aug 20 08:36 network_throughput

      crw-rw-rw-.
      1 root root
      1, 3 Aug 20 08:36 null

      crw-----.
      1 root root
      10, 144 Aug 20 08:36 nvram

      drwxr-xr-x.
      2 root root
      100 Aug 20 08:36 ol_ca-virtdoc-oltest1

      crw-r----.
      1 root kmem
      1, 4 Aug 20 08:36 port

      crw-r----.
      1 root root
      108, 0 Aug 20 08:36 port

      crw-rw-rw-.
      1 root tty
      5, 2 Oct 7 08:10 ptmx

      drwxr-xr-x.
      2 root root
      0 Aug 20 08:36 pts

      crw-rw-rw-.
      1 root root
      1, 8 Aug 20 08:36 pts

      crw-rw-rw-.
      1 root root
      1, 8 Aug 20 08:36 random

      drwxr-xr-x.
      2 root root
      60 Nov 4 2019 raw

      lrwxrwxrwx.
      1 root root
      251, 0 Aug 20 08:36 sda

      brw-rw----.
      1 root disk
      8, 0 Aug 20 08:36 sda

      brw-rw----.
      1 root disk
      8, 1 Aug 20 08:36 sda

      brw-rw----.
      1 root disk
      8, 16 Aug 20 08:36 sdb

      brw-rw----.
      1 root disk
      8, 17 Aug 20 08:36 sdb

    <
```

Block devices support random access to data, seeking media for data, and usually allow data to be buffered while it is being written or read. Examples of block devices include hard disks, CD-ROM drives, flash memory, and other addressable memory devices.

Character devices support the streaming of data to or from a device. The data is not usually buffered nor is random access permitted to data on a device. The kernel writes data to or reads data from a character device one byte at a time. Examples of character devices include keyboards, mice, terminals, pseudo-terminals, and tape drives. tty0 and tty1 are character device files that correspond to terminal devices that allow users to log in from serial terminals or terminal emulators.

Pseudo-terminals secondary devices emulate real terminal devices to interact with software. For example, a user might log in to a terminal device such as /dev/tty1, which then uses the pseudo-terminal primary device, /dev/pts/ptmx, to interact with an underlying pseudo-terminal device. The character device files for pseudo-terminal secondary and primary devices are located in the /dev/pts directory, as shown in the following example:

ls -1 /dev/pts

```
total 0
crw--w---. 1 guest tty 136, 0 Mar 17 10:11 0
crw--w---. 1 guest tty 136, 1 Mar 17 10:53 1
crw--w---. 1 guest tty 136, 2 Mar 17 10:11 2
c----. 1 root root 5, 2 Mar 17 08:16 ptmx
```

Some device entries, such as stdin for the standard input, are symbolically linked through the self subdirectory of the proc file system. The pseudo-terminal device file to which they actually point depends on the context of the process.

ls -1 /proc/self/fd/[012]

```
lrwx-----. 1 root root 64 Oct 7 08:23 /proc/self/fd/0 -> /dev/pts/0 lrwx-----. 1 root root 64 Oct 7 08:23 /proc/self/fd/1 -> /dev/pts/0 lrwx-----. 1 root root 64 Oct 7 08:23 /proc/self/fd/2 -> /dev/pts/0
```

Character devices, such as null, random, urandom, and zero are examples of pseudodevices that provide access to virtual functionality implemented in software rather than to physical hardware. /dev/null is a data sink. Data that you write to /dev/null effectively disappears but the write operation succeeds. Reading from /dev/null returns EOF (end-of-file).

/dev/zero is a data source of an unlimited number of zero-value bytes.

/dev/random and /dev/urandom are data sources of streams of pseudo-random bytes. To maintain high-entropy output, /dev/random blocks if its entropy pool does not contains sufficient bits of noise. /dev/urandom does not block and, as a result, the entropy of its output might not be as consistently high as that of /dev/random. However, neither /dev/random nor /dev/urandom are considered to be truly random enough for the purposes of secure cryptography such as military-grade encryption.

You can find out the size of the entropy pool and the entropy value for /dev/random from virtual files under /proc/sys/kernel/random:

```
cat /proc/sys/kernel/random/poolsize
4096
cat /proc/sys/kernel/random/entropy_avail
3467
```

For more information, see the null (4), pts(4), and random(4) manual pages.

About the Udev Device Manager

The Udev device manager dynamically creates or removes device node files at boot time . When creating a device node, <code>udev</code> reads the device's /sys directory for attributes such as the label, serial number, and bus device number.

Udev can use persistent device names to guarantee consistent naming of devices across reboots, regardless of their order of discovery. Persistent device names are especially important when using external storage devices.

The configuration file for udev is /etc/udev/udev.conf, in which you can define the udev_log logging priority, which can be set to err, info and debug. Note that the default value is err.

For more information, see the udev (7) manual page.

About Udev Rules

Udev uses rules files to determine how it identifies devices and creates device names. The udev service (systemd-udevd) reads the rules files at system start-up and stores the rules in memory. If the kernel discovers a new device or an existing device goes offline, the kernel sends an event action (uevent) notification to udev, which matches the in-memory rules against the device attributes in the /sys directory to identify the device.

Multiple rules files exist in different directories. However, you only need to know about /etc/udev/rules.d/*.rules files because these are the only rules files that you can modify. See Modifying Udev Rules.

Udev processes the rules files in lexical order, regardless of which directory they are located. Rules files in /etc/udev/rules.d override rules files of the same name in other locations.

The following rules are extracted from the file /lib/udev/rules.d/50-udev- default.rules and illustrate the syntax of udev rules:

do not edit this file, it will be overwritten on update



```
SUBSYSTEM=="block", SYMLINK{unique}+="block/%M:%m"
SUBSYSTEM!="block", SYMLINK{unique}+="char/%M:%m"
KERNEL=="pty[pqrstuvwxyzabcdef][0123456789abcdef]", GROUP="tty", MODE="0660"
KERNEL=="tty[pqrstuvwxyzabcdef][0123456789abcdef]", GROUP="tty", MODE="0660"
# mem
KERNEL=="null|zero|full|random|urandom", MODE="0666"
KERNEL=="mem|kmem|port|nvram", GROUP="kmem", MODE="0640"
# block
SUBSYSTEM == "block", GROUP = "disk"
# network
                       MODE="0666"
KERNEL=="tun",
KERNEL=="tun",
KERNEL=="rfkill",
                               MODE="0644"
# CPU
KERNEL=="cpu[0-9]*",
                               MODE="0444"
# do not delete static device nodes
ACTION=="remove", NAME=="", TEST=="/lib/udev/devices/%k", \
    OPTIONS+="ignore remove"
ACTION=="remove", NAME=="?*", TEST=="/lib/udev/devices/$name", \
    OPTIONS+="ignore remove"
```

A rule either assigns a value to a key or it tries to find a match for a key by comparing its current value with the specified value. The following table shows the assignment and comparison operators that you can use.

Operator	Description
=	Assign a value to a key, overwriting any previous value.
+=	Assign a value by appending it to the key's current list of values.
:=	Assign a value to a key. This value cannot be changed by any further rules.
==	Match the key's current value against the specified value for equality.
!=	Match the key's current value against the specified value for equality.

You can use the following shell-style pattern-matching characters in values.

Character	Description
?	Matches a single character.
*	Matches any number of characters, including zero.
[]	Matches any single character or character from a range of characters specified within the brackets. For example, tty[sS][0-9] would match ttys7 or ttyS7.

The following table describes commonly used match keys in rules.



Match Key	Description
ACTION	Matches the name of the action that led to an event. For example, ACTION="add" or ACTION="remove".
$ENV\{key\}$	Matches a value for the device property <i>key</i> . For example, ENV{DEVTYPE}=="disk".
KERNEL	Matches the name of the device that is affected by an event. For example, $\texttt{KERNEL}=="dm-*"$ for disk media.
NAME	Matches the name of a device file or network interface. For example, NAME="?*" for any name that consists of one or more characters.
SUBSYSTEM	Matches the subsystem of the device that is affected by an event. For example, SUBSYSTEM=="tty".
TEST	Tests wheter the specified file or path exists; for example, TEST=="/lib/udev/devices/\$name", where \$name is the name of the currently matched device file.

Other match keys include ATTR{filename}, ATTRS{filename}, DEVPATH, DRIVER, DRIVERS, KERNELS, PROGRAM, RESULT, SUBSYSTEMS, and SYMLINK.

The following table describes commonly used assignment keys in rules.

Assignment Key	Description	
ENV{key}	Specifies a value for the device property <i>key</i> , such as GROUP="disk".	
GROUP	Specifies the group for a device file, such as GROUP="disk".	



Assignment Key	Description	
<pre>IMPORT{type}</pre>	Specifies a set of variables for the device property, depending on <i>type</i> :	
	cmdline	
	Import a single property from the boot kernel command line. For simple flags, udev sets the value of the property to 1. For example, IMPORT {cmdline} = "nodmraid".	
	db	
	Interpret the specified value as an index into the device database and import a single property, which must have already been set by an earlier event. For example, IMPORT{db}="DM_UDEV_LOW_PRIORITY_FLAG".	
	file	
	Interpret the specified value as the name of a text file and import its contents, which must be in environmental key format. For example, IMPORT{file}="keyfile".	
	parent	
	Interpret the specified value as a key-name filter and import the stored keys from the database entry for the parent device. For example IMPORT {parent} = "ID_*".	
	program	
	Run the specified value as an external program and imports its result, which must be in environmental key format. For example IMPORT {program}="usb_idexport %p".	
MODE	Specifies the permissions for a device file, such as MODE="0640".	
NAME	Specifies the name of a device file, such as $\mathtt{NAME}=\mathtt{"em1"}$.	
OPTIONS	Specifies rule and device options, such as OPTIONS+="ignore_remove", which means that the device file is not removed if the device is removed.	
OWNER	Specifies the owner for a device file, such as GROUP="root".	
RUN	Specifies a command to be run after the device file has been created, such as RUN+="/usr/bin/eject \$kernel", where \$kernel is the kernel name of the device.	



Assignment Key	Description
SYMLINK	Specifies the name of a symbolic link to a device file, such as SYMLINK+="disk/by-
	uuid/\$env{ID_FS_UUID_ENC}", where \$env{}
	is substituted with the specified device
	property.

Other assignment keys include $ATTR\{key\}$, GOTO, LABEL, RUN, and WAIT_FOR.

The following table describes the string substitutions that are commonly used with the GROUP, MODE, NAME, OWNER, PROGRAM, RUN, and SYMLINK keys.

String Substitution	Description
<pre>\$attr{file} or %s{file}</pre>	Specifies the value of a device attribute from a file under /sys, such as ENV{MATCHADDR}="\$attr{address}".
<pre>\$devpath or %p</pre>	The device path of the device in the sysfs file system under /sys, such as RUN+="keyboard-force-release.sh \$devpath common-volume-keys".
<pre>\$env{key} or %E{key}</pre>	Specifies the value of a device property, such as SYMLINK+="disk/by-id/md-name-\$env{MD_NAME}-part%n".
\$kernel or %k	Specifies the kernel name for the device.
\$major or %M	Specifies the major number of a device, such as IMPORT{program}="udisks-dm-export %M %m".
<pre>\$minor or %m</pre>	Specifies the minor number of a device, such as RUN+="\$env{LVM_SBIN_PATH}/lvm pvscancachemajor \$majorminor \$minor".
\$name	Specifies the device file of the current device, such as TEST=="/lib/udev/devices/\$name".

Udev expands the strings specified for RUN immediately before its program is executed, which is after udev has finished processing all other rules for the device. For the other keys, udev expands the strings while it is processing the rules.

For more information, see the udev (7) manual page.

Querying Udev and Sysfs

You can use the udevadm command to query the udev database and sysfs.

To query the sysfs device path relative to /sys that corresponds to the device file /dev/sda:

udevadm info --query=path --name=/dev/sda

/devices/pci0000:00/0000:00:0d.0/host0/target0:0:0/0:0:0:0/block/sda



To query the symbolic links that point to /dev/sda, use the following command:

udevadm info --query=symlink --name=/dev/sda block/8:0 disk/by-id/ata-VBOX_HARDDISK_VB6ad0115d-356e4c09 disk/by-id/scsi-SATA_VBOX_HARDDISK_VB6ad0115d-356e4c09 disk/by-path/pci-0000:00:0d.0-scsi-0:0:0:0

To query the properties of /dev/sda, use the following command:

```
udevadm info --query=property --name=/dev/sda
UDEV LOG=3
DEVPATH=/devices/pci0000:00/0000:00:0d.0/host0/target0:0:0/0:0:0/block/sda
MAJOR=8
MINOR=0
DEVNAME=/dev/sda
DEVTYPE=disk
SUBSYSTEM=block
ID ATA=1
ID TYPE=disk
ID BUS=ata
ID MODEL=VBOX HARDDISK
ID REVISION=1.0
ID SERIAL=VBOX HARDDISK VB579a85b0-bf6debae
ID SERIAL SHORT=VB579a85b0-bf6debae
ID ATA WRITE CACHE=1
ID ATA WRITE CACHE ENABLED=1
ID ATA FEATURE SET PM=1
ID ATA FEATURE SET PM ENABLED=1
ID ATA SATA=1
ID ATA SATA SIGNAL RATE GEN2=1
ID SCSI COMPAT=SATA VBOX HARDDISK VB579a85b0-bf6debae
ID PATH=pci-0000:00:0d.0-scsi-0:0:0:0
ID PART TABLE TYPE=dos
LVM SBIN PATH=/sbin
UDISKS PRESENTATION NOPOLICY=0
UDISKS PARTITION TABLE=1
UDISKS PARTITION TABLE SCHEME=mbr
UDISKS PARTITION TABLE COUNT=2
UDISKS ATA SMART IS AVAILABLE=0
DEVLINKS=/dev/block/8:0 /dev/disk/by-id/ata-VBOX HARDDISK VB579a85b0-bf6debae ...
```

To query all of the information for /dev/sda, use the following command:

```
udevadm info --query=all --name=/dev/sda
```

```
P: /devices/pci0000:00/0000:00:0d.0/host0/target0:0:0/0:0:0/block/sda
N: sda
W: 37
S: block/8:0
S: disk/by-id/ata-VBOX_HARDDISK_VB579a85b0-bf6debae
S: disk/by-id/scsi-SATA_VBOX_HARDDISK_VB579a85b0-bf6debae
S: disk/by-path/pci-0000:0d.0-scsi-0:0:0:0
E: UDEV_LOG=3
E: DEVPATH=/devices/pci0000:00/0000:0d.0/host0/target0:0:0/0:0:0/block/sda
E: MAJOR=8
E: MINOR=0
E: DEVNAME=/dev/sda
E: DEVTYPE=disk
E: SUBSYSTEM=block
E: ID ATA=1
```



```
E: ID TYPE=disk
E: ID BUS=ata
E: ID MODEL=VBOX HARDDISK
E: ID SERIAL=VBOX HARDDISK VB579a85b0-bf6debae
E: ID SERIAL SHORT=VB579a85b0-bf6debae
E: ID ATA WRITE CACHE=1
E: ID ATA WRITE CACHE ENABLED=1
E: ID_ATA_FEATURE_SET_PM=1
E: ID ATA FEATURE SET PM ENABLED=1
E: ID ATA SATA=1
E: ID ATA SATA SIGNAL RATE GEN2=1
E: ID SCSI COMPAT=SATA VBOX HARDDISK VB579a85b0-bf6debae
E: ID PATH=pci-0000:00:0d.0-scsi-0:0:0:0
E: ID PART TABLE TYPE=dos
E: LVM SBIN PATH=/sbin
E: UDISKS PRESENTATION NOPOLICY=0
E: UDISKS PARTITION TABLE=1
E: UDISKS PARTITION TABLE SCHEME=mbr
E: UDISKS PARTITION TABLE COUNT=2
E: UDISKS ATA SMART IS AVAILABLE=0
E: DEVLINKS=/dev/block/8:0 /dev/disk/by-id/ata-VBOX HARDDISK VB579a85b0-bf6debae ...
```

To display all of the properties of /dev/sda, as well as the parent devices that udev has found in /sys, use the following command:

udevadm info --attribute-walk --name=/dev/sda

```
looking at device '/devices/pci0000:00/0000:00:0d.0/host0/target0:0:0/0:0:0/block/
sda':
   KERNEL=="sda"
    SUBSYSTEM == "block"
   DRIVER==""
   ATTR{range}=="16"
   ATTR{ext_range}=="256"
   ATTR{removable}=="0"
   ATTR{ro}=="0"
   ATTR{size}=="83886080"
   ATTR{alignment offset}=="0"
   ATTR{capability}=="52"
   ATTR(stat)==" 20884 15437 1254282 338919 5743 8644 103994
109005 ...
   ATTR{inflight}=="
                                     0"
 looking at parent device '/devices/pci0000:00/0000:00:0d.0/host0/
target0:0:0/0:0:0:0':
   KERNELS=="0:0:0:0"
   SUBSYSTEMS=="scsi"
   DRIVERS=="sd"
   ATTRS{device blocked}=="0"
   ATTRS { type } == "0"
   ATTRS{scsi level}=="6"
   ATTRS { vendor } == "ATA
   ATTRS{model}=="VBOX HARDDISK
   ATTRS{rev}=="1.0 "
   ATTRS { state } == "running"
   ATTRS{timeout}=="30"
   ATTRS{iocounterbits}=="32"
   ATTRS{iorequest cnt}=="0x6830"
   ATTRS{iodone cnt}=="0x6826"
   ATTRS{ioerr cnt}=="0x3"
```

```
ATTRS {modalias} == "scsi:t-0x00"
   ATTRS{evt media change} == "0"
   ATTRS{dh state} == "detached"
   ATTRS { queue depth } == "31"
   ATTRS{queue ramp up period}=="120000"
   ATTRS{queue_type}=="simple"
  looking at parent device '/devices/pci0000:00/0000:00:0d.0/host0/target0:0:0':
   KERNELS=="target0:0:0"
   SUBSYSTEMS=="scsi"
   DRIVERS==""
  looking at parent device '/devices/pci0000:00/0000:00:0d.0/host0':
   KERNELS=="host0"
    SUBSYSTEMS=="scsi"
   DRIVERS==""
 looking at parent device '/devices/pci0000:00/0000:00:0d.0':
   KERNELS=="0000:00:0d.0"
   SUBSYSTEMS=="pci"
   DRIVERS=="ahci"
   ATTRS{vendor}=="0x8086"
   ATTRS{device}=="0x2829"
   ATTRS{subsystem vendor}=="0x0000"
   ATTRS{subsystem device} == "0x0000"
   ATTRS{class}=="0x010601"
   ATTRS{irq}=="21"
0,00000003"
   ATTRS{local_cpulist}=="0-1"
   ATTRS{modalias}=="pci:v00008086d00002829sv0000000sd0000000bc01sc06i01"
   ATTRS{numa node}=="-1"
   ATTRS{enable}=="1"
   ATTRS{broken parity status}=="0"
   ATTRS{msi bus}==""
   ATTRS{msi irqs}==""
  looking at parent device '/devices/pci0000:00':
   KERNELS=="pci0000:00"
   SUBSYSTEMS==""
   DRIVERS==""
```

The command starts at the device that is specified by its device path and walks the chain of parent devices. For every device that it finds, it displays all of the possible attributes for the device and its parent devices by using the match key format for udev rules.

For more information, see the udevadm(8) manual page.

Modifying Udev Rules

The order in which rules are evaluated is important. Udev processes rules in lexical order. If you want to add your own rules, you need udev to locate and evaluate these rules before the default rules.

The following example illustrates how to implement a udev rules file that adds a symbolic link to the disk device /dev/sdb.

1. Create a rule file under /etc/udev/rules.d with a file name such as 10-local.rules that udev will read before any other rules file.

The following rule in 10-local.rules creates the symbolic link /dev/my_disk, which points to /dev/sdb:

```
KERNEL=="sdb", ACTION=="add", SYMLINK="my disk"
```

Listing the device files in /dev shows that udev has not yet applied the rule:

ls /dev/sd* /dev/my_disk

```
ls: cannot access /dev/my_disk: No such file or directory /dev/sda /dev/sda1 /dev/sda2 /dev/sdb
```

2. To simulate how udev applies its rules to create a device, you can use the udevadm test command with the device path of sdb listed under the /sys/class/block hierarchy, for example:

udevadm test /sys/class/block/sdb

```
calling: test
version ...
This program is for debugging only, it does not run any program
specified by a RUN key. It may show incorrect results, because
some values may be different, or not available at a simulation run.
...
LINK 'my_disk' /etc/udev/rules.d/10-local.rules:1
...
creating link '/dev/my_disk' to '/dev/sdb'
creating symlink '/dev/my_disk' to 'sdb
...
ACTION=add
DEVLINKS=/dev/disk/by-id/ata-VBOX_HARDDISK_VB186e4ce2-f80f170d
   /dev/disk/by-uuid/a7dc508d-5bcc-4112-b96e-f40b19e369fe
   /dev/my_disk
...
```

3. Restart the systemd-udevd service:

```
sudo systemctl restart systemd-udevd
```

After udev processes the rules files, the symbolic link /dev/my disk has been added:

```
ls -F /dev/sd* /dev/my_disk
/dev/my disk@ /dev/sda /dev/sda1 /dev/sda2 /dev/sdb
```

4. (Optional) To undo the changes, remove /etc/udev/rules.d/10-local.rules and /dev/my disk, then run systemctl restart systemd-udevd again.



5

Managing Kernel Modules

This chapter describes how to load, unload, and modify the behavior of kernel modules.

About Kernel Modules

The boot loader loads the kernel into memory. You can add new code to the kernel by including the source files in the kernel source tree and recompiling the kernel. Kernel modules provide device drivers that enable the kernel to access new hardware, support different file system types, and extend its functionality in other ways. The modules can be dynamically loaded and unloaded on demand. To avoid wasting memory on unused device drivers, Oracle Linux supports loadable kernel modules (LKMs), which enable a system to run with only the device drivers and kernel code that are required to be loaded into memory.

Kernel modules can be signed to protect the system from running malicious code at boot time. When UEFI Secure Boot is enabled, only kernel modules that contain the correct signature information can be loaded. See Oracle Linux: Working With UEFI Secure Boot for more information.

Listing Information About Loaded Modules

The lsmod command lists the modules that are currently loaded into the kernel:

```
1 smod
Module
                 Size Used by
udp diag
                16384 0
ib core
               311296 0
tcp diag
                16384 0
inet diag
                24576 2 tcp_diag,udp_diag
nfsv3
                 49152 0
                16384 1 nfsv3
nfs_acl
24576 0
                 20480 2 dm region_hash,dm_mirror
dm log
. . .
```

The output shows the module name, the amount of memory it uses, the number of processes using the module and the names of other modules on which it depends. The module ${\tt dm_log}$, for instance, depends on the ${\tt dm_region_hash}$ and ${\tt dm_mirror}$ modules. The example also shows that two processes are currently using all three modules.

Display detailed information about a module by using the modinfo command:

modinfo ahci

```
filename: /lib/modules/5.4.17-2136.306.1.3.el8uek.x86_64/kernel/drivers/ata/
ahci.ko.xz
version: 3.0
license: GPL
description: AHCI SATA low-level driver
author: Jeff Garzik
```



srcversion: 3F4E4F52FD2D5F8BBD5F972 alias: pci:v*d*sv*sd*bc01sc06i01*

alias: pci:v00001C44d00008000sv*sd*bc*sc*i*

depends: libahci, libata

retpoline: intree: Y name: ahci

vermagic: 5.4.17-2136.306.1.3.el8uek.x86_64 SMP mod_unload modversions

PKCS#7

sig_id:
signer: Oracle CA Server

sig_key: 22:07:CB:47:59:F3:50:A0:A2:FA:24:CE:B4:00:53:4E:C5:1D:C6:2A

sig_hashalgo: sha512

signature: 2F:AE:AF:6D:56:92:69:C4:77:AB:E1:3D:41:09:AF:A6:FC:1D:3B:A2:

9C:23:79:6F:17:82:D5:A3:9B:61:64:32:72:9B:98:C9:8C:89:73:FB: A4:86:4F:B5:7D:DF:84:8E:05:26:4F:22:CB:02:41:38:7B:7C:CB:C2:

9F:FD:94:8F:35:9B:2A:89:3E:E1:17:40:49:79:30:8B:92:4D:3A:9A:

F4:C7:82:8D:26:BE:6D:FB:71:C6:E5:FD

marvell enable:Marvell SATA via AHCI (1 = enabled) (int) parm: mobile lpm policy: Default LPM policy for mobile chipsets (int) parm:

The output may include the following information:

filename

Absolute path of the kernel object file.

version

Version number of the module. Note that the version number may not be updated for patched modules and may be missing or removed in newer kernels.

license

License information for the module.

description

Short description of the module.

author

Author credit for the module.

srcversion

Hash of the source code used to create the module.

alias

Internal alias names for the module.

depends

Comma-separated list of any modules on which this module depends.

retpoline

A flag indicating that the module is built with support for a mitigation against the Spectre security vulnerability.

intree

A flag indicating that the module is built from the kernel in-tree source and is not tainted.



vermagic

Kernel version that was used to compile the module, which is checked against the current kernel when the module is loaded.

sig id

The method used to store signing keys that may have been used to sign a module for Secure Boot. Usually PKCS#7

signer

The name of the signing key used to sign a module for Secure Boot.

sig key

The signature key identifier for the key used to sign the module.

sig hashalgo

The algorithm used to generate the signature hash for a signed module.

signature

The signature data for a signed module.

parm

Module parameters and descriptions.

Modules are loaded into the kernel from kernel object files (/lib/modules/ kernel_version/kernel/*ko*). To display the absolute path of a kernel object file, specify the -n option, for example:

modinfo -n parport

```
/lib/modules/5.4.17-2136.306.1.3.el8uek.x86 64/kernel/drivers/parport/parport.ko.xz
```

For more information, see the lsmod(5) and modinfo(8) manual pages.

Loading and Unloading Modules

The modprobe command loads kernel modules, for example:

sudo modprobe nfssudo 1smod | grep nfs

```
    nfs
    266415
    0

    lockd
    66530
    1 nfs

    fscache
    41704
    1 nfs

    nfs_acl
    2477
    1 nfs

    auth_rpcgss
    38976
    1 nfs

    sunrpc
    204268
    5 nfs,lockd,nfs_acl,auth_rpcgss
```

Include the $\neg \triangledown$ (verbose) option to show whether any additional modules are loaded to resolve dependencies.

sudo modprobe -v nfs

```
insmod /lib/modules/4.18.0-80.el8.x86_64/kernel/net/sunrpc/auth_gss/auth_rpcgss.ko
insmod /lib/modules/4.18.0-80.el8.x86_64/kernel/fs/nfs_common/nfs_acl.ko
insmod /lib/modules/4.18.0-80.el8.x86_64/kernel/fs/fscache/fscache.ko
...
```





The modprobe command does not reload modules that are already loaded. You must first unload a module before you can load it again.

Use the -r option to unload kernel modules:

```
sudo modprobe -rv nfs
```

```
\label{lib/modules/4.18.0-80.el8.x86_64/kernel/fs/nfs/nfs.kormmod /lib/modules/4.18.0-80.el8.x86_64/kernel/fs/lockd/lockd.kormmod /lib/modules/4.18.0-80.el8.x86_64/kernel/fs/fscache/fscache.ko...
```

Modules are unloaded in reverse order in which they were first loaded. Modules are not unloaded if a process or another loaded module requires them.

For more information, see the modprobe(8) and modules.dep(5) manual pages.

About Module Parameters

To modify a module's behavior, specify parameters for the module in the modprobe command:

```
sudo modprobe module name parameter=value ...
```

Separate multiple parameter/value pairs with spaces. Array values are represented by a comma-separated list, for example:

```
sudo modprobe foo arrayparm=1,2,3,4
```

Alternatively, change the values of some parameters for loaded modules and built-in drivers by writing the new value to a file under <code>/sys/module/module_name/parameters</code>, for example:

```
echo 0 | sudo tee /sys/module/ahci/parameters/skip host reset
```

Configuration files (/etc/modprobe.d/*.conf) specify module options, create module aliases, and override the usual behavior of modprobe for modules with special requirements. The /etc/modprobe.conf file that was used with earlier versions of modprobe is also valid if it exists. Entries in the /etc/modprobe.conf and /etc/modprobe.d/*.conf files use the same syntax.

The following are commonly used commands in modprobe configuration files:

alias

Creates an alternate name for a module. The alias can include shell wildcards. To create an alias for the sd-mod module:

```
alias block-major-8-* sd mod
```

blacklist

Ignore a module's internal alias that is displayed by the modinfo command. This command is typically used if the associated hardware is not required, if two or more



modules both support the same devices, or if a module invalidly claims to support a device. For example, blacklist the alias for the frame-buffer driver cirrusfb:

```
blacklist cirrusfb
```

The /etc/modprobe.d/blacklist.conf file prevents hotplug scripts from loading a module, usually so that a different driver binds the module instead, regardless of which driver happens to be probed first. If it does not already exist, you will need to create it.

install

Runs a shell command instead of loading a module into the kernel. For example, load the module <code>snd-emu10k1-synth</code> instead of <code>snd-emu10k1</code>:

```
install snd-emu10k1 /sbin/modprobe --ignore-install snd-emu10k1 && /sbin/modprobe snd-emu10k1-synth
```

options

Defines options for a module,. To define, for instance, the nohwarypt and gos options for the b43 module:

```
options b43 nohwcrypt=1 qos=0
```

remove

Runs a shell command instead of unloading a module. To unmount /proc/fs/nfsd before unloading the nfsd module:

```
remove nfsd { /bin/umount /proc/fs/nfsd > /dev/null 2>&1 || :; } ;
/sbin/modprobe -r --first-time --ignore-remove nfsd
```

For more information, see the modprobe.conf (5) manual page.

Specifying Modules To Be Loaded at Boot Time

The system loads most modules automatically at boot time. If necessary, you can specify additional modules to be loaded by creating a configuration file for the module in the /etc/modules-load.d directory. The file name should have the extension <code>.conf</code>.

For example to force the <code>bnxt_en.conf</code> to load at boot time, run the following command:

```
echo bnxt_en | sudo tee /etc/modules-load.d/bnxt_en.conf
```

Changes to the /etc/modules-load.d directory persist across reboots.

Preventing Modules From Loading at Boot Time

You can prevent modules from loading at boot time by adding a deny rule in a configuration file in the /etc/modprobe.d directory and then rebuilding the initial ramdisk used to load the kernel at boot time.

1. Create a configuration file to prevent the module from loading. For example:

```
sudo tee /etc/modprobe.d/bnxt_en-deny.conf <<'EOF'
#DENY bnxt_en
blacklist bnxt_en
install bnxt_en /bin/false
EOF</pre>
```

2. Rebuild the initial ramdisk image:



```
sudo dracut -f -v
```

3. Reboot the system for the changes to take effect:

sudo reboot



WARNING:

Disabling modules can have unintended consequences and can prevent a system from booting properly or from being fully functional after boot. It is good practice to create a backup ramdisk image before making changes like this and to take care that your configuration is correct.

About Weak Update Modules

External modules, such as drivers that are installed by using a driver update disk or that are installed from an independent package, are usually installed in the <code>/lib/modules/kernel-version/extra</code> directory. Modules that are stored in this directory are given preference over any matching modules that are included with the kernel when attempting to load them. When installed, these external drivers and modules can override existing kernel modules for the purpose of resolving hardware issues. For each subsequent kernel update, it is critical that these external modules be made available to each compatible kernel so that potential boot issues resulting from driver incompatibilities with the affected hardware can be avoided.

Because the requirement to load the external module with each compatible kernel update is system critical, a mechanism exists to allow external modules to be loaded as weak update modules for compatible kernels.

You make weak update modules available by creating symbolic links to compatible modules in the /lib/modules/kernel-version/weak-updates directory. The package manager handles this process automatically when it detects driver modules that are installed in the /lib/modules/kernel-version/extra directories for any compatible kernels.

For example, if a newer kernel is compatible with a module that was installed for the previous kernel, an external module (such as kmod-kvdo) is automatically added as a symbolic link in the weak-updates directory as part of the installation process, as shown in the following command output:

The existence of the symbolic link allows the external module to loaded for subsequent kernel updates.

In most cases, weak updates make sense and ensure that no extra work is required to carry an external module through subsequent kernel updates. Any potential driver-



related boot issues after kernel upgrades are prevented, thus allowing for more predictable running of a system and its hardware.

In certain cases, you may wish to remove weak update modules in place of a newer kernel, for example, in the case where an issue with a shipped driver has been resolved in a newer kernel. In this case, you might prefer to use the new driver rather than the external module that you installed as part of a driver update.

To remove weak update modules, use the following command:

```
rm -rf /lib/modules/4.18.0-80.el8.x86 64/weak-updates/kmod-kvdo/
```

Running this command manually removes the symbolic links for each kernel.

Alternatively, you can use the weak-modules command, which safely removes the specified weak update module for all of the compatible kernels or the command removes all of the weak update modules for the current kernel. You can use the weak-modules command similarly to add weak update modules.

You can also use the weak-modules command with the dry-run option to test the results without actually running the command, for example:

```
weak-modules --remove-kernel --dry-run --verbose
rm -rf kmod-kvdo
```

For more information about external driver modules and driver update disks, see Oracle Linux 9: Installing Oracle Linux.



6

Configuring Huge Pages

In Oracle Linux, physical memory is managed in fixed-size blocks called pages. In x86_64 architecture, the default size of each page is 4 KB.

The kernel stores virtual to physical address mappings for the pages in a data structure known as the page table. However, page table lookups are resource-intensive, so the most recently used addresses are cached in the CPU's Translation Lookaside Buffer (TLB) for faster retrieval. When the CPU needs to fulfil a request for an address-mapping, it first searches its TLB cache. If the address is found in the TLB cache, it is known as a TLB hit. However, if the requested address-mapping is not found, it is known as a TLB miss, and the system has to carry out a resource-intensive lookup on the page table to retrieve the address information.

The default page size of 4 KB has proved to be suitable for most applications. However, for applications that work with large amounts of memory, the number of 4 KB pages required can be very large and can lead to a high number of TLB misses and a performance overhead. Therefore, Oracle Linux provides huge page features, so that applications requiring more memory can have their requirements fulfilled with fewer pages.

Available Huge Page Features

The following huge page features are supported in Oracle Linux:

HugeTLB Pages

HugeTLB pages are also called static huge pages.

The HugeTLB pages feature enables you to reserve pools of huge pages, each of a specified quantity, for each huge page size supported by the kernel. The currently supported huge page size options on x86_64 architecture are 2 MB and 1 GB.

For more information about configuring HugeTLB pages, see Configuring HugeTLB Pages



Best Practices:

- Make requests to the kernel for static huge pages as close to boot time as possible, as that is when there is least memory fragmentation.
- Huge pages can reduce the amount of memory available to the system.
 Therefore, when requesting a reserved huge page pool, ensure the pool is not oversized and that the system's access to memory is not impacted.

Transparent HugePages

The Transparent HugePages (THP) feature is enabled by default in Oracle Linux. With THP, the kernel automatically assigns huge pages to processes so there is no need to manually reserve the pages. Currently THP only supports 2 MB pages on x86 architecture.

THP can run in the following modes:

- system-wide (default): The kernel attempts to assign huge pages to processes that use large contiguous virtual memory areas whenever it is possible to do so.
- per-process: The kernel only assigns huge pages to application processes that explicitly request huge pages through the madvise() system call.

For more information about configuring THP, see Configuring Transparent HugePages

Configuring HugeTLB Pages

You can configure HugeTLB pages by using the following types of parameters:

- Kernel boot parameters
- File-based configuration parameters

The following sections discuss the parameters in greater detail.

Kernel Boot Parameters for HugeTLB Pages

The kernel boot options enable you to specify values such as the size and the number of pages to be reserved in the kernel's pool. Using the kernel boot parameters is the most reliable method of requesting huge pages.

The following table describes the kernel boot parameters available for HugeTLB page setup.

Table 6-1 The Kernel Boot Command-Line Parameters for Requesting HugeTLB Pages

Parameters	Purpose	Value Options Currently Supported on x86_64 Architecture
default_hugepagesz	Defines the default size of persistent huge pages configured in the kernel at boot time.	2M (default), 1G



Table 6-1 (Cont.) The Kernel Boot Command-Line Parameters for Requesting HugeTLB Pages

Parameters	Purpose	Value Options Currently Supported on x86_64 Architecture
hugepagesz and hugepages	Size parameter hugepagesz is used in conjunction with quantity parameter hugepages to reserve a pool of a specified page size and quantity. For example, to request a pool of 1500 pages of size 2 MB, the command-line options would be as follows:	Hugepagesz: 2M, 1G hugepages: 0 or greater
	hugepagesz=2M hugepages=1500	
	When multiple huge page sizes are supported, the "hugepagesz= <size>hugepages=<qty>" pair can be specified multiple times, once for each supported page size. For example, you can use the following command-line options to request one pool of four pages of 1 GB size, and a second pool of 1500 pages of 2 MB size: hugepagesz=1G hugepages=4 hugepagesz=2M hugepages=1500</qty></size>	

Note:

In a NUMA system, pages reserved with kernel command-line options, as shown in the previous table, are divided equally between the NUMA nodes. If the requirement is to have a different number of pages supported on each node, one possible method is to use the file-based HugeTLB parameters located in the sysfs file system. See File-Based Configuration Parameters for HugeTLB Pages and Requesting HugeTLB Pages Using NUMA Node-Specific Parameters Early in the Boot Process for more information about this.

File-Based Configuration Parameters for HugeTLB Pages

The file-based configuration parameters provide runtime access to the configuration settings.

Note:

In addition to accessing the settings at runtime, you can also initialize the parameters early in the boot process, for instance by creating a start-up bash script or by setting the parameters up in a local rc init script.

There can be multiple instances of each file-based parameter on a system. For example, on a system that supports both 2 MB and 1 GB HugeTLB page sizes, there are multiple instances of nr_hugepages, the parameter for the number of pages in a pool, including the following:

- File /sys/kernel/mm/hugepages/hugepages-2048kB/nr_hugepages for the number of pages in the pool of 2 MB pages.
- File /sys/kernel/mm/hugepages/hugepages-1048576kB/nr_hugepages for the number of pages in the pool of 1 GB pages.

The following table outlines some of the commonly used HugeTLB configuration parameters as well as the multiple file instances that you might find for each parameter.



Table 6-2 Commonly Used File-Based HugeTLB Parameters

Parameter	Purpose	File Paths for Different Instances
nr_hugepages	Each instance of nr_hugepages defines the current number of huge pages in the pool associated with that instance. Can be modified at runtime. Example command: echo 20 sudo tee / proc/sys/vm/nr_hugepages Default value is 0.	The file path formats for different instances of nr_hugepages are as follows: File location: / proc/sys/vm/ nr_hugepages (present on all systems). File location: /sys/ kernel/mm/hugepages/ hugepages- <size>kB/ nr_hugepages (present on systems that support more than one huge page size). File location: /sys/ devices/system/node/ node {0,1,2n}/ hugepages hugepages- <size>kB/nr_hugepages (present on NUMA systems only). Note: You should use the NUMA node-specific path format if you need to request different quantities of pages of different sizes to be supported on specific NUMA nodes. If you use any other path format (for example, /proc/sys/vm/ nr_hugepages) to request HugeTLB pages, the pages will be divided equally between the NUMA nodes.</size></size>



Table 6-2 (Cont.) Commonly Used File-Based HugeTLB Parameters

Parameter	Purpose	File Paths for Different Instances
nr_overcommit_hugepages	 Each instance of nr_overcommit_hugepage s defines the additional number of huge pages, above the quantity specified by nr_hugepages, that can be created by the system at runtime through overcommitting memory. As these additional huge pages become unused, they are then freed and returned to the kernel's normal page pool. Example command: echo 20 sudo tee / proc/sys/vm/nr_overcommit_hugepages 	The file path formats for different instances of nr_overcommit_hugepages are as follows: • File location / proc/sys/vm/ nr_overcommit_hugepage s (present on all systems). • File location: /sys/ kernel/mm/hugepages/ hugepages- <size>kB/ nr_overcommit_hugepage s (present on systems that support more than one huge page size). Note: The nr_overcommit_hugepages parameter is not defined at the individual node level, so there is no node-specific file for this setting.</size>
free_hugepages	Read-only parameter. Each instance of free_hugepages returns the number of huge pages in its associated page pool that have yet to be allocated.	The file path formats for different instances of free_hugepages are as follows: • File location: /sys/ kernel/mm/hugepages/ hugepages- <size>kB/ free_hugepages (present on systems that support more than one huge page size). • File location: /sys/ devices/system/node/ node {0,1,2n}/ hugepages/hugepages-<size>kB/ free_hugepages (present on NUMA systems only).</size></size>



Table 6-2 (Cont.) Commonly Used File-Based HugeTLB Parameters

Parameter	Purpose	File Paths for Different Instances
surplus_hugepages	Read-only parameter. Each instance of surplus_hugepages returns the number of huge pages that have been overcommitted from its associated page pool.	The file path formats for different instances of surplus_hugepages are as follows: • File location: /sys/ kernel/mm/hugepages/ hugepages surplus_hugepages (present on systems that support more than one huge page size). • File location: /sys/ devices/system/node/ node {0,1,2n}/ hugepages/hugepages <pre></pre>

In the following table, each column outlines a file branch under which different instances of the HugeTLB parameters are stored.

Table 6-3 The File Path Locations used for HugeTLB Parameters

/proc/sys/vm	/sys/kernel/mm/hugepages/	/sys/devices/system/node/				
All systems that support static huge pages contain HugeTLB parameter files under / proc/sys/vm.	Systems that support multiple size pools contain HugeTLB parameter files in size-specific folders under /sys/	Only NUMA systems contain HugeTLB parameter files under /sys/devices/system/ node/.				
Note:	kernel/mm/hugepages/.					
On many systems, including many Oracle database servers, the procfs file system is the main parameter-set used. The sysctl parameter vm.nr_hugepages,that is commonly initialized in scripts that request huge pages, also writes to the procfs file /proc/sys/vm/nr_hugepages						



/proc/sys/vm /sys/kernel/mm/hugepages/ /sys/devices/system/node/ The following are example The following are example The following are example folders under branch / folders under branch/sys/ folders under branch /sys/ proc/sys/vm: kernel/mm/hugepages/: devices/system/node: - hugepages-2048kB . . . - free_hugepages - nr_hugepages - nr_hugepages node0 nr_overcommit_hugepages -hugepages nr_overcommit_hugepages hugepages-2048kB - surplus hugepages free hugepages - hugepages-1048576kB nr hugepages free hugepages - nr_hugepages surplus_hugepages hugepages-1048576kB nr overcommit hugepages surplus_hugepages free hugepages nr hugepages surplus hugepages - node1 -hugepages hugepages-2048kB free_hugepages nr hugepages surplus hugepages hugepages-1048576kB free hugepages nr hugepages surplus hugepages

Table 6-3 (Cont.) The File Path Locations used for HugeTLB Parameters



Configuring HugeTLB Pages at Boot Time

The precise way to request huge pages at boot time depends upon the system's requirements. The following example procedures provide possible starting points.

Requesting HugeTLB Pages by Using Kernel Parameters at Boot Time

The following procedure demonstrates how to use kernel command-line options to specify two pools of HugeTLB pages, as well as a default page size, on a system that supports multiple huge page sizes. For the purposes of the example, the following are requested:

- A default page size of 1 GB.
- One pool with four HugeTLB pages of 1 GB size.
- One pool of 1500 HugeTLB pages of 2 MB size.

Before beginning the following procedure, ensure that you have the administrative privileges required for all of the steps.

1. Specify 1 GB size for kernel boot parameter default_hugepagesz, as well as 2 pairs of "hugepagesz=<Size_num>G hugepages=Qty_num" parameters for the two huge page pools, by appending the following line to the kernel command-line options in /etc/default/ grub:

default hugepagesz=1G hugepagesz=1G hugepagesz=2M hugepagesz=2M hugepages=1500

- **2.** Regenerate the GRUB2 configuration file:
 - **a.** If your system uses BIOS firmware, run the following command:

```
sudo grub2-mkconfig -o /boot/grub2/grub.cfg
```

b. If your system uses UEFI framework, run the following command:

```
sudo grub2-mkconfig -o /boot/efi/EFI/redhat/grub.cfg
```

3. The next time the system boots, the two huge page pools will be requested.

Requesting HugeTLB Pages Using NUMA Node-Specific Parameters Early in the Boot Process

Huge Pages requested by using the kernel boot-time parameters, as shown in the previous example, are divided equally between the NUMA nodes.

However, you might need to request a different number of huge pages to be supported on specific nodes. One way of achieving this is to set the configuration values in the node-specific file path defined as $\sys/\devices/\system/node/node\{0,1,2...n\}/\hugepages/\hugepages-<SIZE>kB/.$

The following procedure describes one way of reserving 299 pages of 2 MB size on node 0, and 300 pages of 2 MB size on node 1 on a NUMA system.

Before beginning the following procedure, ensure that you have the administrative privileges required for all of the steps.

 Create a script file called hugetlb-reserve-pages.sh in the /usr/lib/systemd/ directory and add the following content:



2. Make the script executable:

sudo chmod +x /usr/lib/systemd/hugetlb-reserve-pages.sh

3. Create a service file called hugetlb-gigantic-pages.service in the /usr/lib/system/directory and add the following content to it:

```
Description=HugeTLB Gigantic Pages Reservation
DefaultDependencies=no
Before=dev-hugepages.mount
ConditionPathExists=/sys/devices/system/node

[Service]
Type=oneshot
RemainAfterExit=yes
ExecStart=/usr/lib/systemd/hugetlb-reserve-pages.sh

[Install]
WantedBy=sysinit.target
```

4. Enable the service file to ensure that it is called during early boot:

```
sudo systemctl enable hugetlb-gigantic-pages
```

Configuring HugeTLB at Runtime

There may be occasions when you need make a request for huge pages at runtime.

Configuring HugeTLB Pages for a Specific NUMA Node at Runtime

The following example shows how to request 20 HugeTLB pages of size 2048 kB for node2 at runtime.

Before starting the procedure you must ensure you have the required administrative privileges required for all the steps.

1. Run the numastat command to show memory statistics relating to the NUMA nodes:

```
numastat -cm | egrep 'Node|Huge'| grep -v AnonHugePages
```

	Node	0	Node	1	Node	2	Node	3	Total	add
HugePages_Total		0		0		0		0	0	
HugePages_Free		0		0		0		0	0	
HugePages Surp		0		0		0		0	0	

2. Add the required number of huge pages of a specified size to the node of your choice, for example 20 pages of 2MB size on node 2:

```
echo 20 | sudo tee /sys/devices/system/node/node2/hugepages/hugepages-2048kB/ nr_hugepages
```

3. Run the numastat command again to ensure the request was successful and that the requested memory (in our example 20 x 2 MB pages = 40 MB) has been added HugePages Total for node2:

```
        numastat -cm | egrep 'Node|Huge'| grep -v AnonHugePages

        Node 0 Node 1 Node 2 Node 3 Total

        HugePages_Total 0 0 40 0 40

        HugePages_Free 0 0 40 0 40

        HugePages_Surp 0 0 0 0 0 0
```

Configuring Transparent HugePages

The Transparent HugePages (THP) feature is enabled by default in Oracle Linux. However, you may still need to access and configure THP according to your system's needs.

The following sections look at various THP parameters and examples of how they can be configured.

Parameters Used to Configure Transparent HugePages

The following table describes some of the parameter settings that can be used when configuring Transparent HugePages (THP).



Table 6-4 Commonly Used THP Parameters

Parameter	File Location	Value Options
enabled	/sys/kernel/mm/ transparent_hugepage/ enabled	The enabled parameter determines if THP is enabled or not, and if enabled, which mode it is running in. The following list explains the available options: • always (default): THP is enabled in system— wide mode. In this setting, the kernel attempts to assign huge pages to processes using large contiguous virtual memory areas whenever it is possible to do so.
		 madvise: THP is enabled in per- process mode. In this setting the kernel only assigns huge pages to application processes that explicitly request huge pages through the madvise() system call. disabled: THP is disabled.



Table 6-4 (Cont.) Commonly Used THP Parameters

Parameter	File Location	Value Options
defrag	/sys/kernel/mm/ transparent_hugepage/ defrag	The defrag parameter determines how aggressively an application should reclaim pages and defrag memory when THP is unavailable. The following list explains the available options: always: An application requesting THP will stall on allocation failure and directly reclaim pages and compact memory in an attempt to obtain a THP immediately. defer: An application does not stall but continues using small pages. The application requests the kernel daemons kswapd and kcompactd to reclaim pages and compact memory so that THP is available later. defer+madvise: Regions using the madvise (MADV_HUGEPAGE) call will stall on allocation failure and directly reclaim pages and compact memory in an attempt to obtain a THP immediately However, all other regions will requests the kernel daemons kswapd and kcompactd to reclaim pages and compact memory in an attempt to obtain a THP immediately However, all other regions will requests the kernel daemons kswapd and kcompactd to reclaim pages and compact memory so that THP is available later. madvise (default): Regions using the madvise (MADV HUGEPAG
		E) call will stall on allocation failure and directly reclaim pages and compact memory in an attempt to obtain



Configuring Transparent HugePages at Runtime

The following sections show you examples of how you can configure THP at runtime by accessing the THP parameters in the sysfs virtual file system.

Retrieving the Current Status of Transparent HugePages

To see the current setting of THP you can read the /sys/kernel/mm/ transparent hugepage/enabled parameter as shown in the following code sample:

```
sudo cat /sys/kernel/mm/transparent_hugepage/enabled
[always] madvise never
```

The value in the output that is enclosed in square brackets represents the current setting.

Changing the Current Status of Transparent HugePages

To change the current status of THP you need to write the setting of your choice to/sys/kernel/mm/transparent_hugepage/enabled. The following example shows you how to set the status to always:

1. Check the current status of THP by reading the enabled parameter.

```
sudo cat /sys/kernel/mm/transparent_hugepage/enabled
always madvise [never]
```

The value in the output that is enclosed in square brackets represents the current setting.

Set THP mode to always.

```
echo always | sudo tee /sys/kernel/mm/transparent hugepage/enabled
```

3. Confirm the change has been successful by reading the enabled parameter.

```
sudo cat /sys/kernel/mm/transparent_hugepage/enabled
[always] madvise never
```

The value in the output enclosed in square brackets represents the current setting.



Note:

Virtual file systems like sysfs provide a file system interface to items that are not necessarily stored as files on disk. The sysfs files therefore do not always interact with file commands in exactly the same way that regular physical files on disk would. For instance, in the previous example, the echo command used does not overwrite $/sys/kernel/mm/transparent_hugepage/enabled$, as it would if used with a regular file, but instead simply changes the selected option:

```
sudo cat /sys/kernel/mm/transparent_hugepage/enabled
    always madvise [never]
echo always | sudo tee /sys/kernel/mm/transparent_hugepage/enabled
    always
sudo cat /sys/kernel/mm/transparent_hugepage/enabled
    [always] madvise never
```

Changing the defrag Setting of Transparent HugePages

To change the THP defrag setting you need to write the setting of your choice to $/sys/kernel/mm/transparent_hugepage/defrag$



The optimum <code>defrag</code> setting varies from system to system. Reclaiming pages and memory compaction can increase the number of THP pages available, but the reclaiming and compaction also take up CPU time, so you need to find the correct balance for your system.

The following example shows you how to set the defrag setting to madvise.

1. Check the current value of the defrag parameter:

```
sudo cat /sys/kernel/mm/transparent_hugepage/defrag
[always] defer defer+madvise madvise never
```

The value in the output that is enclosed in square brackets represents the current setting.

2. Set the /sys/kernel/mm/transparent hugepage/defrag parameter to madvise:

```
echo madvise | sudo tee /sys/kernel/mm/transparent_hugepage/defrag
```

3. Confirm the change has worked by reading the defrag parameter.

```
sudo cat /sys/kernel/mm/transparent_hugepage/defrag
always defer defer+madvise [madvise] never
```

The value in the output enclosed in square brackets represents the current setting.



7

Managing Resources

This chapter describes how to manage the use of resources in an Oracle Linux system.

About Control Groups

Control groups (cgroups) are a collection of resources, such as CPU, memory, network, and so on, which have collective settings on how these resources are used by applications and processes. Cgroups constitute a functionality in Oracle Linux that enables you to further manage the use of resources by setting limits, prioritizing, isolating, or allocating these resources. Thus, you are able to control resource use on a granular level to obtain a more efficient system performance. Control groups are particularly important in system configurations that host multiple virtual machines, Kubernetes clusters, and so on, whose applications compete over resource use.

Oracle Linux supports two types of control groups:

Control groups version 1 (cgroups v1)

These groups provide a per-resource controller hierarchy. Each resource, such as CPU, memory, I/O, and so on, has its own control group hierarchy. A disadvantage of this group is the difficulty of establishing proper coordination of resource use among groups that might belong to different process hierarchies.

Control groups version 2 (cgroups v2)

These groups provide a single control group hierarchy against which all resource controllers are mounted. In this hierarchy, you can obtain better proper coordination of resource uses across different resource controllers. This version is an improvement over cgroups v1 whose over flexibility prevented proper coordination of resource use among the system consumers.

Both versions are present in Oracle Linux. However, by default, the <code>cgroups v2</code> functionality is enabled and mounted on Oracle Linux 9 systems.

For more information about control groups of both versions, see the cgroups (7) and sysfs (5) manual pages.

About Kernel Resource Controllers

Control groups manage resource use through *kernel resource controllers*. A kernel resource controller represents a single resource, such as CPU time, memory, network bandwidth, or disk I/O.

To identify the currently mounted resource controllers in the system, check the contents of the /procs/cgroups file, for example:

sudo less /proc/cgroups
#subsys_name hierarchy num_cgroups enabled
cpuset 0 103 1
cpu 0 103 1

cpuacct	0	103	1	
blkio	0	103	1	
memory	0	103	1	
devices	0	103	1	
freezer	0	103	1	
net_cls	0	103	1	
perf_eve	ent	0	103	1
net_prid		0	103	1
hugetlb	0	103	1	
pids	0	103	1	
rdma	0	103	1	
misc	0	103	1	

For a detailed explanation of the kernel resource controllers of both <code>cgroups v1</code> and <code>cgroups v2</code>, see the <code>cgroups(7)</code> manual page.

About the Control Group File System

In Oracle Linux, the <code>cgroup</code> functionality is mounted as a file system in /sys/fs/cgroup. This directory is also called the root control group. The contents of the directory differ depending on which <code>cgroup</code> version is mounted on the system. For <code>cgroups v2</code>, the directory contents are as follows:

ls /sys/fs/cgroup

```
cgroup.controllerscpuset.mems.effectivememory.statcgroup.max.depthcpu.statmisc.capacitycgroup.max.descendantsdev-hugepages.mountsys-fs-fuse-connections.mountcgroup.procsdev-mqueue.mountsys-kernel-config.mountcgroup.statinit.scopesys-kernel-debug.mountcgroup.subtree_controlio.pressuresys-kernel-tracing.mountcgroup.threadsio.statsystem.slicecpu.pressurememory.numa_statuser.slicecpuset.cpus.effectivememory.pressure
```

To create a new control group, create a child group or subdirectory in the root control group.

```
sudo mkdir /sys/fs/cgroup/MyGroup
```

The child group is automatically populated with resource controllers that you have enabled for the new group. For sample procedures that create child groups where you can implement resource management for an application, see Setting CPU Weight to Regulate Distribution of CPU Time and Setting CPU Bandwidth to Regulate Distribution of CPU Time.

To destroy a child group, ensure that the child group itself does not contain additional child groups, then remove the directory from the root control group.

```
sudo rm -rf /sys/fs/cgroup/MyGroup
```

About Control Groups and systemd

Control groups can be used by the systemd system and service manager for resource management. Systemd uses these groups to organize units and services that consume resources. For more information about systemd, see About the systemd Service Manager.



Systemd supports different unit types, three of which are for resource control purposes:

- **Service**: A process or a group of processes whose settings are based on a unit configuration file. Services encompass specified processes in a "collection" so that systemd can start or stop the processes as one set. Service names follow the format name.service.
- **Scope**: A group of externally created processes, such as user sessions, containers, virtual machines, and so on. Similar to services, scopes encapsulate these created processes and are started or stopped by the arbitrary processes and then registered by systemd at runtime. Scope names follow the format name.scope.
- Slice: A group of hierarchically organized units in which services and scopes are located. Thus, slices themselves do not contain processes. Rather, the scopes and services in a slice define the processes. Every name of a slice unit corresponds to the path to a location in the hierarchy. Root slices, typically user.slice for all user-based processes and system.slice for system-based processes, are automatically created in the hierarchy. Parent slices exist immediately below the root slice and follow the format parent-name.slice. These root slices can then have sub-slices on multiple levels.

The service, the scope, and the slice units directly map to objects in the control group hierarchy. When these units are activated, they map directly to control group paths that are built from the unit names. To display the mapping between the systemd resource unit types and control groups, type:

```
sudo systemd-cqls
Working directory /sys/fs/cgroup:
-user.slice (#1243)
 → trusted.invocation id: 50ce3909b2644f919ee420adc39edb4b
  ⊢user-1001.slice (#4167)
   → trusted.invocation id: 02e80a960d4549a7a9c69ce0fb546c26
    ⊢session-2.scope (#4405)
      -2417 sshd: oracle [priv]
       -2430 sshd: oracle@pts/0
       -2431 -bash
       -2689 sudo systemd-cgls
       -2691 systemd-cqls
      └_2692 less
    └user@984.service ... (#3827)
     → trusted.delegate: 1
     → trusted.invocation id: 09b47ce9f3124239b75814114353f3f2
      └init.scope (#3861)
        -2058 /usr/lib/systemd/systemd --user
        └_2099 (sd-pam)
 -init.scope (#19)
  └─1 /usr/lib/systemd/systemd --switched-root --system --deserialize 17
 -system.slice (#53)
   -chronyd.service (#2467)
   → trusted.invocation id: c0f77aaa9c7844e6bef6a6898ae4dd56
    └1358 /usr/sbin/chronyd -F 2
   -auditd.service (#2331)
   → trusted.invocation id: 756808add6a348609316c9e8c1801846
   └1310 /sbin/auditd
   -tuned.service (#3079)
   → trusted.invocation id: 2c358135fc46464d862b05550338d4f4
   └1415 /usr/bin/python3 -Es /usr/sbin/tuned -l -P
   -systemd-journald.service (#1651)
```



For an example of how to use systemd commands such as systemct1 to manage resources, see Controlling Access to System Resources. For further technical details, see the systemct1(1), systemd-cgls(1), and systemd.resource-control(5) manual pages.

About Resource Distribution Models

The following distribution models provide you ways of implementing control or regulation in distributing resources for use by cgroups v2:

Weights

In this model, the weights of all the control groups are totaled. Each group receives a fraction of the resource based on the ratio of the group's weight against the total weight.

Consider 10 control groups, each with a weight of 100 for a combined total of 1000. In this case, each group can use a tenth of a specified resource.

Weight is typically used to distribute stateless resources. To apply this resource, the $\mbox{CPUWeight}$ option is used.

Limits

A control group can use the configured amount of a resource. However, you can also overcommit resources. Therefore, the sum of the subgroups limits can exceed the limit of the parent group.

To implement this distribution model, the MemoryMax option is used.

Protections

In this model, a group is assigned a "protected boundary". Provided that the group's resource usage remains within the protected amount, the kernel cannot deprive the group of the use of the resource in favor of other groups that are competing for the same resource. In this model, an overcommitment of resources is permitted. To implement this model, the MemoryLow option is used.

Allocations

In this model, a specific absolute amount is allocated for the use of finite type of resources, such as real-time budget.

Using cgroups v2 to Manage Resources for Applications

This section describes how to use $\mathtt{cgroups}\ \mathtt{v2}$ to manage the resource use of applications and processes that are running in the system. The procedures specifically focus on regulating CPU use by these applications so that CPU consumption becomes more efficient.



To prevent applications from overuse of CPU time, you can use <code>cgroups v2</code> to set limits so that the applications' use of the resource is better regulated. In the sections that follow, CPU resource control is implemented by employing two methods:

- Setting CPU bandwidth
- Setting CPU weight

Enabling cgroups v2

At boot time, Oracle Linux 9 mounts cgroups v2 by default.

1. Verify that cgroups v2 is enabled and mounted on the system.

```
sudo mount -1 | grep cgroup

cgroup2 on /sys/fs/cgroup type cgroup2
(rw,nosuid,nodev,noexec,relatime,seclabel,nsdelegate,memory_recursiveprot)
```

2. Optionally, check the contents of /sys/fs/cgroup directory, which is also called the root control group.

```
ll /sys/fs/cgroup/
```

For cgroups v2, the files in the directory should have prefixes to their file names, for example, cgroup.*, cpu.*, memory.*, and so on. See About the Control Group File System.

Preparing the Control Group for Distribution of CPU Time

1. Verify that in the root control group, the cpu and cpuset controllers are available in the /sys/fs/cgroup/cgroup.controllers file.

```
sudo cat /sys/fs/cgroup/cgroup.controllers
cpusetcpu io memory hugetlb pids rdma misc
```

2. Add the CPU controllers to the cgroup.subtree control file.

By default, only the memory and pids controllers are in the file. To add the CPU controllers, type:

```
echo "+cpu" | sudo tee /sys/fs/cgroup/cgroup.subtree_control
echo "+cpuset" | sudo tee /sys/fs/cgroup/cgroup.subtree control
```

3. Optionally, verify that the CPU controllers have been properly added.

```
sudo cat /sys/fs/cgroup/cgroup.subtree_control
cpuset cpu memory pids
```

4. Create a child group under the root control group to become the new control group for managing CPU resources on applications.

```
sudo mkdir /sys/fs/cgroup/MyGroup
```

5. Optionally, list the contents of the new subdirectory or child group.

```
11 /sys/fs/cgroup/MyGroup
```

```
-r-r--. 1 root root 0 Jun 1 10:33 cgroup.controllers
-r-r--. 1 root root 0 Jun 1 10:33 cgroup.events
-rw-r-r-. 1 root root 0 Jun 1 10:33 cgroup.freeze
-rw-r-r-. 1 root root 0 Jun 1 10:33 cgroup.max.depth
```



```
-rw-r-r-. 1 root root 0 Jun 1 10:33 cgroup.max.descendants
-rw-r-r-. 1 root root 0 Jun 1 10:33 cgroup.procs
-r-r--. 1 root root 0 Jun 1 10:33 cgroup.stat
-rw-r-r-. 1 root root 0 Jun 1 10:33 cgroup.subtree control
-rw-r-r-. 1 root root 0 Jun 1 10:33 cpuset.cpus
-r-r-r-. 1 root root 0 Jun 1 10:33 cpuset.cpus.effective
-rw-r-r-. 1 root root 0 Jun 1 10:33 cpuset.cpus.partition
-rw-r-r-. 1 root root 0 Jun 1 10:33 cpuset.mems
-r-r-r-. 1 root root 0 Jun 1 10:33 cpuset.mems.effective
-r-r-r-. 1 root root 0 Jun 1 10:33 cpu.stat
-rw-r-r-. 1 root root 0 Jun 1 10:33 cpu.weight
-rw-r-r-. 1 root root 0 Jun 1 10:33 cpu.weight.nice
-r-r-r-. 1 root root 0 Jun 1 10:33 memory.events.local
-rw-r-r-. 1 root root 0 Jun 1 10:33 memory.high
-rw-r-r-. 1 root root 0 Jun 1 10:33 memory.low
-r-r-r-. 1 root root 0 Jun 1 10:33 pids.current
-r-r-r-. 1 root root 0 Jun 1 10:33 pids.events
-rw-r-r-. 1 root root 0 Jun 1 10:33 pids.max
```

Based on the CPU controllers that you added to <code>/sys/fs/cgroup/cgroup.subtree_control</code>, the contents of the <code>MyGroup</code> that are inherited from the root control group are now more limited. Thus, only <code>cpuset.*</code>, <code>cpu.*</code>, <code>memory.*</code>, and <code>pids.*</code> files are in the <code>MyGroup</code> directory.

6. Enable the CPU-related controllers in MyGroup's cgroup. subtre control files.

```
echo "+cpu" | sudo tee /sys/fs/cgroup/MyGroup/cgroup.subtree_control
echo "+cpuset" | sudo tee /sys/fs/cgroup/MyGroup/cgroup.subtree control
```

 Optionally, verify that the CPU controllers are enabled for child groups under MyGroup.

```
sudo cat /sys/fs/cgroup/MyGroup/cgroup.subtree_control
cpuset cpu
```

Setting CPU Bandwidth to Regulate Distribution of CPU Time

This procedure is based on the following assumptions:

• The application that is consuming CPU resources heavily is shalsum, as shown in the following sample output of the top command:

```
sudo top
 PID USER PR NI VIRT RES SHR S %CPU %MEM
                                             TIME+
COMMAND
 34500 root 20 0 18720 1756 1468 R 99.0 0.0 0:31.09
sha1sum
 34501 root 20 0 18720 1772 1480 R 99.0
                                       0.0
                                           0:30.54
sha1sum
   1 root 20 0 109724 17196 11032 S 0.0
                                       0.1
                                           0:03.28
svstemd
   2 root 20 0
                               0 S 0.0 0.0 0:00.00
kthreadd
    3 root 0 -20
                     0
                          0 0 I 0.0 0.0 0:00.00
rcu gp
```



```
4 root 0 -20 0 0 0 I 0.0 0.0 0:00.00 rcu_par_gp ...
```

In the sample output, the shalsum processes have PIDs 34500 and 34501, respectively.

The current system has multiple CPUs.

To display the number of CPUs in the system, you can use the following command:

```
sudo cat /sys/fs/cgroup/cpuset.cpus.effective
0-1
```

The sample output indicates a dual-core system.

Important:

As a prerequisite to the following procedure, you must complete the preparations of cgroup-v2 as described in Preparing the Control Group for Distribution of CPU Time. If you skipped those preparations, you cannot complete this procedure.

1. Create a tasks directory in the MyGroup subdirectory.

```
sudo mkdir /sys/fs/cgroup/MyGroup/tasks
```

This directory defines a child group with files that relate only to cpu and cpuset controllers.

2. Optionally, list the contents of the new subdirectory.

```
11 /sys/fs/cgroup/MyGroup/tasks
```

```
11 /sys/fs/cgroup/Example/tasks
-r-r-r-. 1 root root 0 Jun 1 11:45 cgroup.controllers
-r-r-r-. 1 root root 0 Jun 1 11:45 cgroup.events
-rw-r-r-. 1 root root 0 Jun 1 11:45 cgroup.freeze
-rw-r-r-. 1 root root 0 Jun 1 11:45 cgroup.max.depth
-rw-r-r-. 1 root root 0 Jun 1 11:45 cgroup.max.descendants
-rw-r-r-. 1 root root 0 Jun 1 11:45 cgroup.procs
-r-r-r-. 1 root root 0 Jun 1 11:45 cgroup.stat
-rw-r-r-. 1 root root 0 Jun 1 11:45 cgroup.subtree control
-rw-r-r-. 1 root root 0 Jun 1 11:45 cgroup.threads
-rw-r-r-. 1 root root 0 Jun 1 11:45 cgroup.type
-rw-r-r-. 1 root root 0 Jun 1 11:45 cpu.max
-rw-r-r-. 1 root root 0 Jun 1 11:45 cpu.pressure
-rw-r-r-. 1 root root 0 Jun 1 11:45 cpuset.cpus
-r-r--. 1 root root 0 Jun 1 11:45 cpuset.cpus.effective
-rw-r-r-. 1 root root 0 Jun 1 11:45 cpuset.cpus.partition
-rw-r-r-. 1 root root 0 Jun 1 11:45 cpuset.mems
-r-r-r-. 1 root root 0 Jun 1 11:45 cpuset.mems.effective
-r-r--. 1 root root 0 Jun 1 11:45 cpu.stat
-rw-r-r-. 1 root root 0 Jun 1 11:45 cpu.weight
-rw-r-r-. 1 root root 0 Jun 1 11:45 cpu.weight.nice
-rw-r-r-. 1 root root 0 Jun 1 11:45 io.pressure
-rw-r-r-. 1 root root 0 Jun 1 11:45 memory.pressure
```

3. In the tasks directory, set the processes that you want to regulate for CPU time to use the same CPU.

```
echo "1" | sudo tee /sys/fs/cgroup/MyGroup/tasks/cpuset.cpus
```

4. Optionally, verify that the processes run on the same CPU.

```
cat /sys/fs/cgroup/MyGroup/tasks/cpuset.cpus
```

5. Configure the CPU bandwidth to set restrictions within the MyGroup/tasks child control group.

```
echo "200000 1000000"n | sudo tee /sys/fs/cgroup/MyGroup/tasks/cpu.max
```

In the command, the value 20000 represents the quota of time in microseconds that is allowed for all processes collectively in a child group to run during a specified period. That period, in turn, is defined by the value 1000000. Specifically, the processes in the /sys/fs/cgroup/MyGroup/tasks group can run on the CPU for only 0.2 seconds, or one fifth, of every second.

If the quota is exhausted by the control group within the defined period, then the processes are suspended until the next period.

6. Optionally, verify the time quotas.

top

```
sudo cat /sys/fs/cgroup/MyGroup/tasks/cpu.max
200000 1000000
```

7. Add the PIDs of the applications to the child group.

```
echo "34500" | sudo tee /sys/fs/cgroup/MyGroup/tasks/cgroup.procs echo "34501" | sudo tee /sys/fs/cgroup/MyGroup/tasks/cgroup.procs
```

8. Verify that the applications are running in the specified control group.

```
sudo cat /proc/34500/cgroup /proc/34501/cgroup
0::/MyGroup/tasks
0::/MyGroup/tasks
```

9. Check the current CPU consumption after you have set the CPU bandwidth.

PID USER PR NI VIRT RES SHR S %CPU %MEM TIME+ COMMAND 34500 root 20 0 18720 1756 1468 R 10.0 0.0 37:36.13 sha1sum 34501 root 20 0 18720 1772 1480 R 10.0 0.0 37:41.22 sha1sum 1 root 20 0 186192 13940 9500 S 0.0 0.4 0:01.60 systemd 2 root 20 0 0 0 0 S 0.0 0.0 0:00.01 kthreadd 0 -20 0 -20 3 root 0 I 0 0 0.0 0.0 0:00.00 rcu gp 4 root 0 0 0 I 0.0 0.0 0:00.00 rcu par gp . . .

Because the MyGroup/tasks group is limited to a total of 20% of CPU use, then each shalsum process is now limited to 10% of CPU time.



Setting CPU Weight to Regulate Distribution of CPU Time

This procedure is based on the following assumptions:

• The application that is consuming CPU resources heavily is shalsum, as shown in the following sample output of the top command:

```
sudo top
  PID USER PR NI VIRT RES SHR S %CPU %MEM TIME+ COMMAND
 33301 root
            20 0 18720 1756 1468 R 99.0 0.0 0:31.09 sha1sum
 33302 root
            20 0 18720 1772 1480 R 99.0 0.0 0:30.54 sha1sum
 33303 root 1 root
            20 0 18720 1772 1480 R 99.0 0.0 0:30.54 shalsum
                                           0.1 0:03.28
            20 0 109724 17196 11032 S 0.0
systemd
    2 root
            20 0
                                  0 S 0.0
                                           0.0 0:00.00
kthreadd
    3 root
                                 0.0
             0 -20
                            0
                                           0.0 0:00.00
             0 -20 0
    4 root
                           0
                                0 I 0.0 0.0 0:00.00
rcu par gp
```

In the output, the shalsum processes have PIDs 33301, 33302, and 33303, respectively.

The current system has multiple CPUs.

To display the number of CPUs in the system, you can use the following command:

```
sudo cat /sys/fs/cgroup/cpuset.cpus.effective
0-1
```

The sample output indicates a dual-core system.

Important:

As a prerequisite to the following procedure, you must complete the preparations of cgroup-v2 as described in Preparing the Control Group for Distribution of CPU Time. If you skipped those preparations, you cannot complete this procedure.

1. Create 3 child groups in the MyGroup subdirectory.

```
sudo mkdir /sys/fs/cgroup/MyGroup/g1
sudo mkdir /sys/fs/cgroup/MyGroup/g2
sudo mkdir /sys/fs/cgroup/MyGroup/g3
```

2. Configure the CPU weight for each child group.

```
echo "150" | sudo tee /sys/fs/cgroup/MyGroup/g1/cpu.weight echo "100" | sudo tee /sys/fs/cgroup/MyGroup/g2/cpu.weight echo "50" | sudo tee /sys/fs/cgroup/MyGroup/g3/cpu.weight
```

3. Apply the application PIDs to their corresponding child groups.



```
echo "33301" | sudo tee /sys/fs/cgroup/Example/g1/cgroup.procs echo "33302" | sudo tee /sys/fs/cgroup/Example/g2/cgroup.procs echo "33303" | sudo /sys/fs/cgroup/Example/g3/cgroup.procs
```

These commands set the desired applications to become members of the $MyGroup/g^*/control$ groups. The CPU time for each shalsum process depends on the CPU time distribution as configured for each group.

The weights of the g1, g2, and g3 groups that have running processes are summed up at the level of MyGroup, which is the parent control group.

With this configuration, when all processes run at the same time, the kernel allocates to each of the shalsum processes the proportionate CPU time based on their respective cgroup's cpu.weight file, as follows:

Child group	cpu.weight setting	Percent of CPU time allocation			
g1	150	~50% (150/300)			
g2	100	~33% (100/300)			
g3	50	~16% (50/300)			

If one child group has no running processes, then the CPU time allocation for running processes is recalculated based on the total weight of the remaining child groups with running processes. For example, if the g2 child group does not have any running processes, then the total weight becomes 200, which is the weight of g1+g3. In this case, the CPU time for g1 becomes 150/200 (~75%) and for g3, 50/200 (~25%)

4. Check that the applications are running in the specified control groups.

```
sudo cat /proc/33301/cgroup /proc/33302/cgroup /proc/33303/cgroup
0::/MyGroup/g1
0::/MyGroup/a2
```

0::/MyGroup/g2
0::/MyGroup/g3

5. Check the current CPU consumption after you have set the CPU weights.

СОР											
					222	2110	2	0.0011	0.4534	TT14T :	
COMMAND	USER	PR	NI	VIRT	RES	SHR	S	%CPU	%MEM	TIME+	
33301 shalsum	root	20	0	18720	1748	1460	R	49.5	0.0	415:05.87	
33302	root	20	0	18720	1756	1464	R	32.9	0.0	412:58.33	
sha1sum 33303	root	20	0	18720	1860	1568	R	16.3	0.0	411:03.12	
sha1sum		0.0		41.6600	00540	15006	_		0 5	0 10 00	
/60	root	20	0	416620	28540	15296	S	0.3	0.7	0:10.23	tuned
1	root	20	0	186328	14108	9484	S	0.0	0.4	0:02.00	
systemd											
2	root	20	0	0	0	0	S	0.0	0.0	0:00.01	
kthread											



Using cgroups v2 to Manage Resources for Users

The previous sample procedures describe how to manage applications' use of system resources. You can also manage resource use by directly implementing resource filters to users who log in to the system. Run Control Groups Version 2 on Oracle Linux is a tutorial that provides examples on how to control users' use of system resources. Additionally, the tutorial offers a lab environment that enables you to perform steps in real time to regulate resource consumption by users.

